



Topological inference from measures

Mickaël Buchet

► To cite this version:

Mickaël Buchet. Topological inference from measures. Computational Geometry [cs.CG]. Université Paris Sud - Paris XI, 2014. English. NNT : 2014PA112367 . tel-01108521

HAL Id: tel-01108521

<https://inria.hal.science/tel-01108521>

Submitted on 22 Jan 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Comprendre le monde,
construire l'avenir®



UNIVERSITÉ PARIS-SUD
ÉCOLE DOCTORALE 427
INFORMATIQUE DE PARIS-SUD

THÈSE

présentée pour obtenir

LE GRADE DE DOCTEUR EN SCIENCES

DE L'UNIVERSITÉ PARIS SUD

Spécialité : Informatique

par

Mickaël Buchet

Sujet :

Topological inference from measures

Date de soutenance : 01 décembre 2014

Composition du jury :

Directeur de thèse :	M. Frédéric Chazal	Directeur de recherche (INRIA - Saclay Île-de-France)
Co-directeur de thèse :	M. Steve Oudot	Chargé de recherche (INRIA - Saclay Île-de-France)
Rapporteurs :	Mme. Dominique Attali	Directrice de recherche (CNRS - Gipsa-lab - Grenoble)
	M. Xavier Goaoc	Professeur (Université de Marne-la-Vallée)
Examineurs :	M. Tamal Dey	Professeur (The Ohio State University)
	M. Marc Schoenauer	Directeur de recherche (INRIA - Saclay Île-de-France)

Abstract

In this thesis, we study the problem of analysing topological structure in point cloud data. One widely used tool in this domain is persistent homology. By processing the data at all scales, it does not rely on a particular choice of scale, which is one of the main challenge faced in this area. Moreover, its stability properties provide a natural connection between discrete data and an underlying continuous structure. Finally, it can be combined with other tools, like the distance to a measure, which allows to handle noise that are unbounded. The main caveat of this approach is its high complexity.

In this thesis, we will introduce topological data analysis and persistent homology, then show how to use approximation to reduce the computational complexity. We provide an approximation scheme to the distance to a measure and a sparsifying method of weighted Vietoris-Rips complexes in order to approximate persistence diagrams with practical complexity. We detail the specific properties of these constructions.

Persistent homology was previously shown to be of use for scalar field analysis. We provide a way to combine it with the distance to a measure in order to handle a wider class of noise, especially data with unbounded errors. Finally, we discuss interesting opportunities opened by these results to study data where parts are missing or erroneous.

Contents

Abstract	iii
Acknowledgments	ix
1 Introduction	1
1.1 Persistent homology	2
1.2 The problem of outliers	4
1.3 The distance to a measure	5
1.4 Contributions	6
1.5 Organisation	8
2 Preliminaries	9
2.1 Simplicial and singular homologies	9
2.1.1 Simplicial complexes	9
2.1.2 Chain complexes	10
2.1.3 Simplicial homology	11
2.1.4 Simplicial maps	13
2.1.5 Singular homology	13
2.2 Persistence	14
2.2.1 Persistence diagrams	14
2.2.2 Bottleneck distance	16
2.2.3 Persistence module interleaving	18
2.2.4 Computation of persistence diagrams	20
2.3 Geometry	21
2.3.1 Metric spaces	21
2.3.2 Compact sets	22
2.3.3 Riemannian geometry	23
2.3.4 Lower bound on the volume of Riemannian balls	24
3 Distance to a measure	27
3.1 Measures	27
3.2 Wasserstein distances	28
3.3 Distance to a measure	29
3.4 Stability	31

3.5	Discriminating results	34
3.5.1	Correspondence between measure and distance to a measure	34
3.5.2	Recovering a measure from its distance	35
3.5.3	Application to the real line	37
3.5.4	Higher dimensional reconstruction for finite support measures	39
3.6	Relation to higher order Voronoi and power diagrams	40
3.6.1	k^{th} -order Voronoi diagrams	41
3.6.2	Power diagrams	44
4	Distance to a measure approximation	47
4.1	Barycentric decomposition	47
4.2	Random sampling	49
4.3	Witnessed k -distance	51
4.4	Power distance with compact support	54
4.4.1	Stability of power distances	55
4.4.2	Approximation of the distance to a measure	57
4.4.3	Restriction to measures with finite support	58
4.4.4	Euclidean case	59
4.5	Application to persistence diagrams approximation	61
4.6	Example of approximations	63
5	Persistence of power distance functions approximation	67
5.1	Weighted Rips filtration	67
5.1.1	Definition	67
5.1.2	Stability	68
5.1.3	Approximation of $d_{\mu,m}$	72
5.2	Weighted Rips induced metric	72
5.3	Sparse weighted Rips	74
5.3.1	Sparse Rips complexes	76
5.3.2	Projection onto Nets	78
5.3.3	Sometimes the projections induce contiguous simplicial maps	79
5.3.4	Sparse filtrations and power distance functions	81
5.4	Experimental illustration	82
6	Specific noise conditions for the distance to a measure	89
6.1	Counter examples for the Wasserstein noise condition	89
6.2	New noise conditions for distances to measures	92
6.3	Relation to other noise conditions	93
6.3.1	Clutter noise	94
6.3.2	Wasserstein noise condition	95
6.3.3	Sampling by empirical measures	96
6.3.4	Discrete results for Hausdorff noise condition	97
6.4	Consequences on persistence diagrams	99

7	Scalar field analysis	103
7.1	Scalar field analysis with bounded noise	103
7.2	Unbounded functional noise	104
7.2.1	Functional noise model	105
7.2.2	Functional denoising	108
7.2.3	Experimental Illustration for functional noise	111
7.3	Scalar field analysis with unbounded geometric noise	115
7.4	Scalar field analysis with both functional and geometric noise	116
8	Regression and incomplete data	121
8.1	Discrepancy is a regression estimator	121
8.1.1	The discrepancy regression	121
8.1.2	Convergence rate for discrepancy	122
8.1.3	Convergence rate for median	123
8.1.4	Relaxing the noise model	124
8.2	Application to incomplete data	126
8.2.1	Algorithm for recovery of incomplete data	126
8.2.2	Illustration on a synthetic example	127
9	Conclusion	129
A	Proof of Theorem 5.12	131
B	Résumé en français	141
B.1	Introduction	141
B.2	Homologie persistante	142
B.3	Le problème du bruit aberrant	145
B.4	La distance à la mesure	146
B.5	Contenu de la thèse	147
B.5.1	Distance à la mesure	147
B.5.2	Structures de données pour la persistance des distances de puissance	149
B.5.3	Nouvelles conditions d'échantillonnage pour la distance à la mesure	150
B.5.4	Analyse de champs scalaires et données incomplètes	151
B.6	Organisation de la thèse	151

Acknowledgments

A dissertation like this one is the result of more than three years of individual work. I would not have succeeded to produce this document and the research it contains without guidance and collaborations. I want to express my gratitude to all those that participated, either with their collaborations, discussions or support.

First, I am grateful to my advisers, Frédéric Chazal and Steve Oudot. Their high standards in mathematical writing sometimes led to pertinent comments I had trouble to cope with. They endured my temper and kept me on the right path. Without them, I would, without doubt, have wandered in research without completing any work.

I would also like to thank Dominique Attali and Xavier Goaoc who took the time to carefully review my thesis and made useful comments on how to improve it. My thanks extend to Marc Schoenauer and Tamal Dey who accepted to be on the committee.

Concerning Tamal Dey, I have multiple reasons to be grateful. My visit to Columbus was made possible by his invitation along with Yusu Wang. I had a very enlightening stay at The Ohio State University and I hope they enjoyed our collaboration as much as I did. During this visit, I was warmly welcomed by their team and I want to specifically thank Fengtao and Issam. You took care of me outside working hours and managed to keep me out of trouble and boredom. I hope to be able to repay your kindness someday.

The geometric team is split between Sophia-Antipolis and Saclay. I had the chance to have worthy collaborations with different researchers on both sites. I thank especially Don Sheehy, Marc Glisse, Jean-Daniel Boissonnat, Olivier Devillers and Mariette Yvinec. In Saclay, there is an invaluable member of the team: our assistant, Christine Biard, does a wonderful job supporting the the team and goes beyond her work, lending an ear when I needed to speak. Fellow graduate students, both inside and outside the team, were comrades to fight the loneliness of the PhD. I will without doubt forget some of them but I appreciated the time spent with Cécile, Thomas, Ruqi, Étienne, as well as the numerous discussions with Amaury, Jonas, Jean-Philippe, Mikaël and Clément.

Outside the academic world, I will not forget Roxane, who was there in the most difficult times of my PhD, as well as my family, parents and siblings, for their occasional comments and continuous support.

And to all those I forgot, I beg your forgiveness.

Palaiseau, November 2014

MICKAËL BUCHET

1 Introduction

Data gathering is now a daily activity, from companies collecting data about their customers or employees, to intelligence agencies and polling organisations. The belief is that information is a resource. For example, some governments want to encourage economic growth through the use of data as shown by the open data initiative [1, 2, 3]. However, data in itself is useless. It is necessary to interpret it and shape it into information.

Interpretation is often done through visualisation. Given a 2 or 3-dimensional object, we, human beings, are able to interpret it. Data is usually given as a point set in some high dimensional space. For example, grey scale images are points in a space whose dimension is the number of pixels. An answer to a poll is a point in a space whose dimension is the number of questions. Such high-dimensional data are impossible to visualise directly and we need to process them before interpretation.

Multiple problems are part of this interpretation. Clustering [41, 59, 78, 95] and segmentation [98] try to separate points into different groups. Reconstruction tries to recover a continuous object from data points, usually under the form of a triangulation [21, 47]. Dimensionality reduction projects data onto a smaller-dimensional space, which can make it easier to visualise or analyse, using the most relevant parameters to describe the data [62, 93].

In this thesis, we consider topological data analysis and more precisely topology inference. We aim to recover structure in the data by inferring the underlying topology. This knowledge can guide us for the resolution of the above problems. If we know the correct number of connected components, then we know the right number of clusters we should obtain in clustering. If we know the intrinsic dimension of the data, then we know the size of the space we need to use for dimensionality reduction. One of the most popular tools for topological data analysis in recent years has been persistent homology, which analyses the data at different scales.

One must not forget that practical data is almost always noisy, either from measurement errors or from imperfect models. Topological data analysis methods need to be robust against noise. Existing algorithms usually work well when the noise is bounded. However, aberrant values are common in data. A faulty sensor or a mistake can create points that have no relation with the rest of the data and are difficult to handle.

Recently, persistent homology has spread to a large variety of fields. One application has been to use persistent homology to define signatures for data. In this setting, persistence provides topological information that discriminates between different classes of phenomena. This

has been used to classify images of different pathologies [4, 38], to analyse electroencephalograms [97] or differentiate between shapes [23]. Moreover it can provide a way to cluster or segment data [33, 85, 90]. The next step is to search for a certain pattern and thus detect and identify features, which has been applied for images [77], subtypes of cancer [84] or cyclic patterns in genome [46].

Persistent homology also provides a way to better understand the structure of objects and visualise it, from the structure of matter in astrophysics [91, 92] to compressed granular media [75], complex networks [72, 86] or dynamical systems [10]. In biology, it can explain protein compressibility [65] and describe root structures [57]. It has also been used to study the propagation of genes coding resistances to antibiotics [58]. The reconstruction itself can be done [35], provide structure for tracking [9] and visualisation for cortical structures [76, 89]. The output of persistent homology is usually a persistence diagram, a structure that does not suit well a statistical setting. For example, we do not know how to define the mean of two diagrams. However, the introduction of persistence landscapes [16] makes possible the use of statistical analysis for topology, as has been done for orthodontic data [64, 70].

1.1 Persistent homology

Given a family of nested topological spaces indexed by a parameter $\alpha \in \mathbb{R}$, $\mathcal{F} = \{F_\alpha\}$, persistent homology studies the evolution of the topology of these spaces as α grows from $-\infty$ to ∞ . In data analysis, the most usual way to build a sequence of topological spaces is to grow balls. Given a point cloud P , it means that we are considering the sub-level sets of the distance to P . Assuming that P is sampling an underlying object K , the hope is that some of the sub-level sets have the same topological type as K . The topology of the sub-level sets is often stable over an interval of values. In this thesis, we only consider parameters defined over subsets of \mathbb{R} .

By topology, we mean the homology. Intuitively, it corresponds to the connected components in dimension 0, the holes or cycles in dimension 1, the cavities in dimension 2 and so on. Consider Figure 1.1, where we have a set of points noisily sampled along the edges of a square S . We want to recover the topology of S which has 1 connected component and 1 cycle. For $\alpha < 0$, the sub-level set of the distance to P is empty. When $\alpha = 0$, the sub-level set is exactly P and therefore has 14 connected components, 1 for each point of P . As α increases and the balls grow, we finally obtain a sub-level set that has the same topology as S . Remark that it is stable for some values α .

The topological information obtained using persistence is usually represented by a persistence diagram. A topological feature, for example a cycle, appears in one of the topological spaces of $F_\alpha \in \mathcal{F}$. α is called the *birth time* of the topological feature. The topological feature then exists in some F_γ , $\alpha < \gamma < \beta$ and no longer exists in F_δ for $\delta > \beta$. β is called the *death time* of the topological feature. Note that one 0-dimensional feature, id est, one connected component does not die and thus has an infinite death time. The persistence diagram of dimension d of \mathcal{F} is the multi-set composed of pairs (x, y) where x is the birth time of d -dimensional feature and y is the death time of the same dimensional feature. Persistence diagrams can be represented either by a multi-set in \mathbb{R}^2 or a barcode. In the first case, every pair (x, y) is represented by a point. In the second case, (x, y) is represented by a bar starting at x and ending at y . Figure 1.2

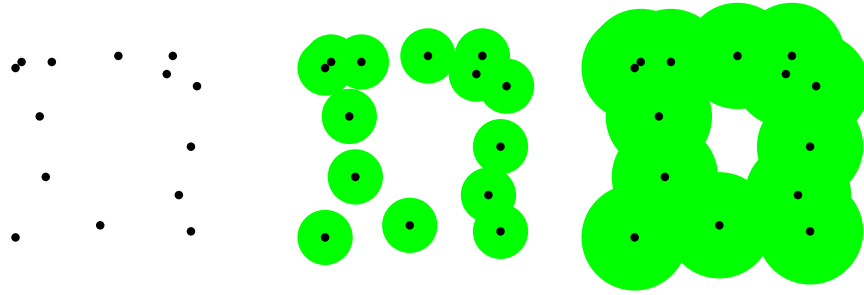


Figure 1.1 – Growing balls for persistence

shows the two representations obtained for topological features of dimension 1 (cycles) in the example of Figure 1.1.

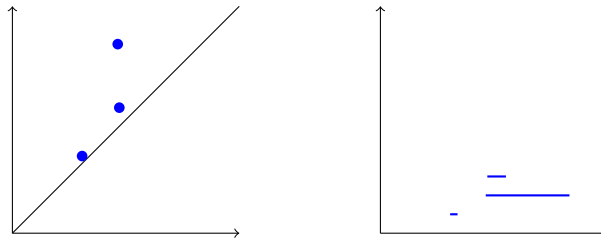


Figure 1.2 – Persistence diagram in dimension 1

The idea of persistence is that the topology features that corresponds to the underlying object K are stable over an interval of values of the parameter. Therefore, they have a longer lifespan than the feature due to noise. In our example, we see that one feature of dimension 1 has longer lifespan than the others. It corresponds to the cycle in S . The two smaller bars are due to small cycles appearing as the balls grow.

Persistence is multi-scale. We consider the whole range of values for α and therefore, we look at the data at all scales. It means that we can detect and recover the topology of objects that have different topology depending on the considered scale. For example, the point cloud of Figure 1.3 samples a spiral rolled over a torus. When we look at it very closely, we have a 1-dimensional object, the spiral. From an intermediate distance, we have the torus, which is a 2-dimensional object. Persistent homology is able to correctly analyse this difference of topology depending on the scale.

A good persistence diagram for inference is a diagram where the ratio between the lifespan of relevant features and irrelevant ones, called *gap*, is large. Persistence diagrams are stable to small variations of the function used to defined the sub-level sets. When the distance to P approximates the distance to K , we obtain a good diagram.

The computation of persistence diagrams does not escape the so-called curse of dimensionality, which means that they work well in small dimensions but not in higher ones due to a blow up in complexity. The technique to compute persistence diagrams of a family \mathcal{F} of union

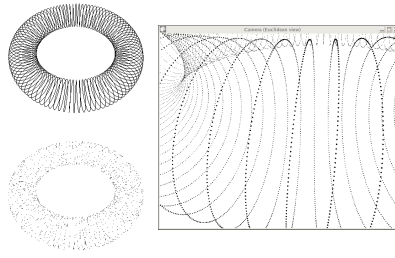


Figure 1.3 – Spiral over a torus

of balls is to build an increasing family of simplicial complexes \mathcal{G} . A simplicial complex is a set of points, edges, triangles, tetrahedra and so on. The family \mathcal{G} approximates the topology of \mathcal{F} . The classical algorithm [57] computing persistent homology has a time complexity of $O(N^3)$, where N is the number of simplices in the maximal simplicial complex in \mathcal{G} . However, if \mathcal{F} is described using n balls in a space of dimension d , we need to build the maximal d -dimensional simplicial complex in \mathcal{G} . Its size is $\binom{n}{d}$. Thus the complexity will be of order $O(n^{3d})$, which make it unusable in practice for high dimensions.

Recent approaches try to make the complexity dependent on the intrinsic dimension instead of the extrinsic one. For example, this has been achieved for Vietoris-Rips complexes [88]. It means that an object of small dimension embedded in a much higher dimensional space can be analysed without paying the complexity cost of the ambient space. This is a way to circumvent the curse of dimensionality without doing dimensionality reduction.

1.2 The problem of outliers

Aberrant values, also called outliers, create problems when computing persistence diagrams. Consider the 1-skeleton of a cube, id est, the set of its edges. In input, we are given a point set that samples the skeleton and contains four outliers located at the centre of four of the cube faces, such that the two empty faces are opposite, as shown in Figure 1.4. These noisy points perturb the persistence diagram significantly and reduce the gap.

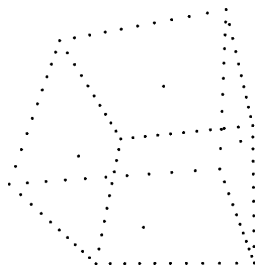


Figure 1.4 – Sampling of a cube skeleton with outliers

We aim at recovering the persistence diagram of the cube skeleton shown in Figure 1.5, id est the persistence diagram of the sub-level sets of the distance to the cube skeleton. The object has a unique connected component appearing at 0 and existing for all non negative values of the parameter α . At the start, we have 5 topological features of dimension 1, or cycles, because

the cube has six faces and one is algebraically the sum of the five others. As the offset grows, the faces are filled and the 1-dimensional features disappear, replaced by a 2-dimensional feature corresponding to the void inside the cube.

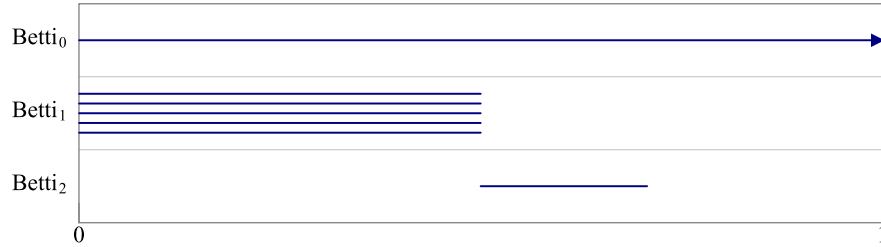


Figure 1.5 – Cube skeleton persistence diagram

The presence of outliers disrupts the persistence diagram. Computing the persistence diagram of the sub-level sets filtration of the distance to the point cloud, we obtain Figure 1.6. Observe that the diagram in dimension 1 has now a smaller gap but we can recover the correct structure. However, in dimension 2, the persistent homology is completely different and the gap is 1, which means that we can not separate signal from noise. We have two topological features. Each corresponds to one half of the cube. When the faces are filled by the growing offsets, a connection simultaneously forms in the middle of the cube, created by the four outliers.

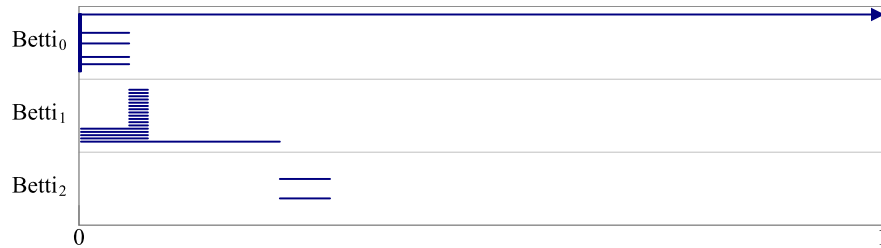


Figure 1.6 – Persistence diagram obtained from the sampling with outliers

Recovering the whole diagram is especially useful when persistent homology is used to derive signatures. Differences in the later part of the diagram can provide interesting information to discriminate objects. However, outliers can completely change the aspect of the diagram.

1.3 The distance to a measure

To handle noise and especially outliers, the idea is to replace the distance to the point cloud P by another function. Such a function must have two properties. It needs to be stable with respect to small variations in the data and its sub-level sets have to be easily computable.

We use the distance to a measure. Given a point set P in a metric space \mathbb{X} with n points and a mass parameter $m = \frac{k}{n}$ where k is an integer, the distance to the empirical measure μ on P is

the function defined over \mathbb{X} by

$$d_{\mu,m}(x) = \sqrt{\frac{1}{k} \sum_{i=1}^k d_{\mathbb{X}}(x, p_i(x))^2}$$

where $p_i(x)$ is the i^{th} nearest neighbour of x in P and $d_{\mathbb{X}}(x, p_i(x))$ is the distance between x and $p_i(x)$. In a more general setting, the distance to a measure μ quantifies the cost of the best transport plan to bring a fraction m of the mass from μ to the point we consider. This function is stable as has been first shown in [25] and can be easily computed in Euclidean spaces [67]. Thus it is a good candidate for topological inference.

Going back to the cube skeleton, the persistence diagram of the distance to the empirical measure and a mass corresponding to 5 points is given in Figure 1.7. The diagram still contains noise but there exists a clear difference between the lifespans of real topological features and those due to noise in dimension 1 and 2.

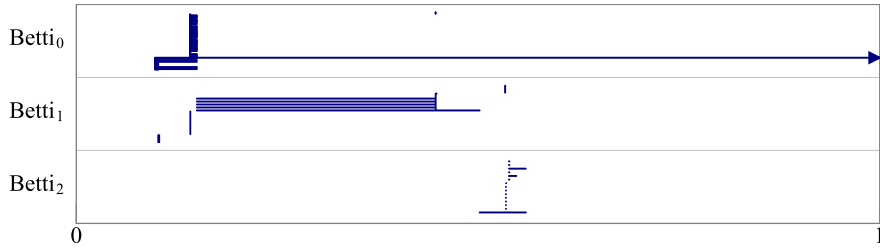


Figure 1.7 – Diagram obtained using the distance to the empirical measure for the cube skeleton

The use of the persistence diagram of the distance to a measure for real data is confronted with a major issue. Outside of Euclidean spaces, the sub-level sets of the distance to a measure are not computable. In Euclidean spaces, the sub-level sets of the distance to a measure are a union of balls [67]. However, the number of balls needed to describe them is the same as the number of non empty cell in the k^{th} -order Voronoi diagram of P , which can be as large as $O\left(n^{\lfloor \frac{d+1}{2} \rfloor} k^{\lceil \frac{d+1}{2} \rceil}\right)$ [39]. Hence, it is necessary to approximate them before being able to use it in practice. An approximation with a linear number of balls was proposed in [67] and lower bounds on the number of balls needed are given in [80]. Unfortunately, these results are limited to Euclidean spaces and do not extend to other metric spaces.

This approximation is not sufficient to handle data in a high-dimensional space. The large size of the simplicial complex used to compute persistence diagrams remains. The result from [88] assumes that for one $F_\alpha \in \mathcal{F}$, all balls have the same radius. The balls in the sub-level sets filtration of $d_{\mu,m}$ do not verify this property and the method of [88] has to be adapted.

1.4 Contributions

This thesis investigates the complexity of computing persistent homology and how we can handle noise with aberrant values. We aim at making homology inference robust to noise,

and tractable for intrinsically low-dimensional data, even if embedded in high dimensional spaces. We propose an approximation method to the persistence diagrams of the distance to a measure and introduce new settings in which it can be used in order to broaden the set of possible applications.

Distance to a measure. The distance to a measure μ was defined in Euclidean spaces in [25]. Its stability is guaranteed if two measures are close with respect to the Wasserstein distance. We straightforwardly extend this stability result to general metric spaces. Once the stability of the distance to a measure is established, a natural question concerns the identifiability of measures using the distances to measures. We provide new results on how we can identify measures knowing only their distance to measure functions.

The second requirement for the use of the distance to a measure in persistent homology is the computability of its sub-level sets. We improve the theoretical guarantees of the first approximation provided in [67]. The sub-level sets are approximated using a linear number of balls with a constant multiplicative approximation factor. We introduce new ways to approximate distance to a measure functions in any metric space with similar guarantees. We provide tight theoretical bounds.

Data structures for persistence of power distances. The result of the approximation of the distance to a measure is a power distance. The computation of persistence diagrams of power distance is thus an important challenge. The Vietoris-Rips filtration classically used to approximate persistence diagrams for distance functions can be adapted to the power distance case by the addition of weights. It gives a new structure called the weighted Rips filtration. We study the stability of this construction and show that it is a metric space.

However, the use of the weighted Rips is not enough to make the computation of persistence tractable. Like its unweighted version, it does not escape the curse of dimensionality. The size of the filtration needed to compute the persistence diagram is exponential in the ambient dimension. We adapt the linear sized approximation of the Vietoris-Rips filtration [88] to the weighted Rips filtration. We thus obtain a filtration whose size is linear in the number of input points and exponential in the intrinsic dimension of the data and such that its persistence diagram approximates the persistence diagram of the power distance.

New noise conditions for the distance to a measure. Given two probability measures μ and ν , the assumptions on the Wasserstein distance between μ and ν to ensure the closeness of the two distances to measures $d_{\mu,m}$ and $d_{\nu,m}$ are not optimal. We propose new conditions such that the distance to a measure can be used for the computation of persistence diagrams. In particular, it can allow a partial recovery of the persistence diagram when a complete approximation is not possible.

Scalar field analysis and incomplete data. Persistence is also known to be used for scalar field analysis [32]. The aim is to study the structure of a real valued function defined over a Riemannian manifold. The previous algorithm was unable to handle outliers, either in the

function values or the position of the points. Using a new estimator for function values, built using the distance to a measure, we adapt the previous pipeline to handle combinations of noise with outliers in the geometry as well as in the function values. Moreover, this estimator can be seen as a regression operator and we provide convergence rates.

This kind of techniques seems promising for the analysis of incomplete data. We provide an elementary algorithm and some examples suggesting further directions of research. The theoretical setting for incomplete data does not exist yet and we only present an illustration of the problem.

1.5 Organisation

The research presented in this thesis is partially the outcome of two collaborations soon to be published.

Chapters 2 and 3 introduce classical notions in topological data analysis and straightforwardly extend the stability results of the distance to a measure to general metric spaces. Section 3.5, however, presents new results on the identifiability of measures using the distance to a measure.

Chapters 4 and 5 are the result of a collaboration with F. Chazal, S. Oudot and D. Sheehy [18]. The first one concerns approximation to the distance to a measure, while the second one studies the approximation of the persistence diagrams of power distance functions. The approximation of the distance to a measure is obtained by an original method which also provide new results for existing methods. The data structure for the computation of persistence diagram is a technical adaptation of a previous work from D. Sheehy [88].

Chapter 7 studies scalar field analysis and is the outcome of a visit to T. Dey and Y. Wang at The Ohio State University. The results, obtained in collaboration with F. Chazal, S. Oudot and F. Fan [17] are an adaptation of [32] to new noise conditions, using the distance to a measure and a new function estimator.

Finally, Chapters 6 and 8 present connected work studying the noise conditions of Chapter 7 and the proprieties of the new estimator. Chapter 8 also presents an opening to the incomplete data problem.

2 Preliminaries

In this chapter, we introduce notions of algebraic topology and geometry used in this thesis. First, we formalise the topological properties we study. They are the simplicial and singular homologies. Then, we define their persistent versions before making some general geometric considerations.

2.1 Simplicial and singular homologies

We introduce two algebraic constructions. First is the simplicial homology, built using simplicial complexes. It provides a computable descriptor of the topology of simplicial complexes. Second is the singular homology, which allows to discuss the homology of more general spaces and can be related to simplicial homology.

2.1.1 Simplicial complexes

The construction of simplicial and singular homologies rely on simplicial complexes. We first introduce abstract simplicial complexes.

Definition 2.1 *Given a set of indices $I \subset \mathbb{N}$, an abstract n -simplex is a set of $n + 1$ distinct elements of I .*

Definition 2.2 *Given a set of indices $I \subset \mathbb{N}$, an abstract simplicial complex is a set S of abstract simplices such that, for any subset $\sigma \in S$, every simplex $\sigma' \subset \sigma$ belongs to S .*

Abstract simplices and complexes can be realised using geometric simplices.

Definition 2.3 *Given an Euclidean space \mathbb{R}^d , a geometric n -simplex σ is the convex hull of a set of n points (v_0, \dots, v_n) in \mathbb{R}^d .*

If the resulting object has dimension n , the n -simplex is said to be *non-degenerate*. Given a n -simplex σ , the $n + 1$ points of which σ is the convex hull are the *vertices* of σ . Consider the $n + 1$ sets of points obtained by removing one point from the vertices of σ . Each of these sets defines a $(n - 1)$ -simplex σ' called a *facet* of σ .

Figure 2.1 shows a 2-simplex and its facets of decreasing dimensions. A 2-simplex is a triangle. It has 3 facets of dimension 1 which are its edges. These edges have each 2 facets of dimension 0 which are the vertices of the triangle.

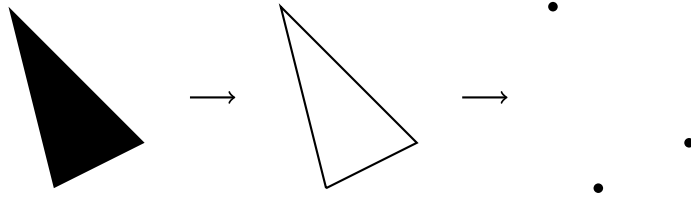


Figure 2.1 – Decomposition of a 2-simplex

Definition 2.4 A set of geometric simplices S is a geometric complex, if for any n -simplex $\sigma \in S$, all facets σ' of σ belong to S and if the intersection of any pair of simplices is a face of each of them.

Every abstract simplicial complex is isomorphic to a geometric simplicial complex [83, Theorem 3.1]. This simplicial complex is called the *geometric realisation* of the abstract complex and is a topological space for the topology induced by inclusion. The combinatorial structure of abstract and geometric complexes is equivalent. We use the geometric complexes to give some intuition about homology.

Simplicial complexes sharing the same 1-skeleton, i.e. the same set of edges or 1-simplices, can be ordered using the inclusion. This ordering is partial but has a greatest element called the *clique complex*.

Definition 2.5 Let P be a set of points and $E \subset P^2$ a set of edges. The clique complex of (P, E) is the maximal simplicial complex for the inclusion among complexes whose 0-simplices are P and 1-simplices are E .

In other words, the clique complex contains all simplices that could be built using the set of edges E . Clique complexes form a family of simplicial complexes with some interesting properties. First, they can be stored efficiently as only the 1-skeleton is needed to describe the whole complex. Secondly, it is relatively easy to prove the contiguity of simplicial maps involving clique complexes using Lemma 2.12 and thus to provide theoretical guarantees for algorithms in Chapter 5.

2.1.2 Chain complexes

Homology is an algebraic construction using a chain complex. Consider a sequence of Abelian groups $\{C_i\}_{i \geq 0}$ and homomorphisms $\{\partial_i\}_{i \geq 0}$ between those groups such as shown in Figure 2.2. A *chain complex* is a sequence where, for all n , $\partial_n \partial_{n+1} = 0$. This condition means that $\text{Im } \partial_{n+1} \subset \text{Ker } \partial_n$ and we can define quotient groups using these spaces.

$$\cdots \longrightarrow C_{n+1} \xrightarrow{\partial_n} C_n \longrightarrow \cdots \longrightarrow C_1 \xrightarrow{\partial_1} C_0 \xrightarrow{\partial_0} 0$$

Figure 2.2 – Chain complex

Definition 2.6 Given a chain complex C , the quotient group $\text{Ker } \partial_n / \text{Im } \partial_{n+1}$ is called the n^{th} homology group of C .

Building a quotient group means that all elements of $\text{Im } \partial_{n+1}$ are equivalent to 0. This implies that two elements are equivalent if and only if their difference is an element of $\text{Im } \partial_{n+1}$. The elements of $\text{Im } \partial_{n+1}$ are called *boundaries* while the elements of $\text{Ker } \partial_n$ are called *cycles*. All homologies use this construction and differ by the choice of the groups $\{C_i\}$ and homomorphisms $\{\partial_i\}$.

2.1.3 Simplicial homology

We consider a simplicial complex X and assume that we have a numbering of the vertices of X . For the purpose of simplicial homology, the order of the vertices in simplices has an importance. It defines a notion of orientation. Given a n -simplex σ , we write $[v_0, \dots, v_i, \dots, v_j, \dots, v_n]$ the set of its vertices. The orientation is reversed when two vertices are exchanged, $-\sigma = [v_0, \dots, v_j, \dots, v_i, \dots, v_n]$.

Let $\Delta_n(X)$ be the free Abelian group whose basis is the n -simplices of X with coefficient in a ring \mathbb{A} . Writing $S_n(X)$ for the set of all n -simplices of X , the elements of $\Delta_n(X)$ are of the form $\sum_{\sigma \in S_n(X)} n_\sigma \sigma$. Given two elements $\sum n_\sigma \sigma$ and $\sum n'_\sigma \sigma$, their sum is given by $\sum (n_\sigma + n'_\sigma) \sigma$.

$$\dots \longrightarrow \Delta_{n+1}(X) \xrightarrow{\partial_n} \Delta_n(X) \longrightarrow \dots \longrightarrow \Delta_1(X) \xrightarrow{\partial_1} \Delta_0(X) \xrightarrow{\partial_0} 0$$

Figure 2.3 – Chain complex over a simplicial complex

The boundary operator is induced by the decomposition of simplices shown in Figure 2.1. In dimension 2, the boundary of a triangle will be its three edges. Signs in the coefficient n_σ intuitively indicate the orientation of the simplices and the boundary operator is defined such that the orientation of all simplices is coherent. Given a n -simplex σ and its vertices $[v_0, \dots, v_n]$, we denote $[v_0, \dots, \hat{v}_i, \dots, v_n]$ the simplex obtained by removing the vertex v_i .

Definition 2.7 Given a simplicial complex X , the boundary homomorphism $\partial_n : \Delta_n(X) \rightarrow \Delta_{n-1}(X)$ is defined on the basis elements by:

$$\partial_n(\sigma) = \sum_i (-1)^i [v_0, \dots, \hat{v}_i, \dots, v_n]$$

The sign part of the definition ensures that the orientation is consistent as shown in Figure 2.4. The definition over the basis set extends to a homomorphism. By convention, $\partial_n = 0$ for $n \leq 0$. Remark that applying ∂_1 to the three edges from the triangle will give 0 because each vertex appears once positively and once negatively. It is consistent with the construction of a chain complex, which requires that $\partial_{n-1} \circ \partial_n = 0$.

Lemma 2.8 The composition $\partial_{n-1} \circ \partial_n : \Delta_n(X) \rightarrow \Delta_{n-2}(X)$ is zero.

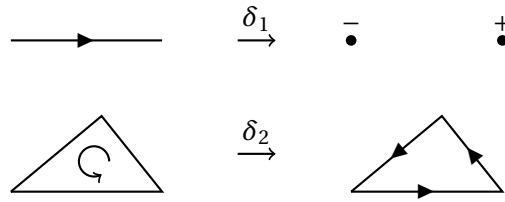


Figure 2.4 – Examples of boundaries

Proof: Let σ be a n -simplex. By definition, $\partial_n(\sigma) = \sum_i (-1)^i [v_0, \dots, \hat{v}_i, \dots, v_n]$ and hence

$$\begin{aligned} \partial_{n-1} \circ \partial_n(\sigma) &= \sum_{j < i} (-1)^i (-1)^j [v_0, \dots, \hat{v}_j, \dots, \hat{v}_i, \dots, v_n] + \sum_{j > i} (-1)^i (-1)^{j-1} [v_0, \dots, \hat{v}_i, \dots, \hat{v}_j, \dots, v_n] \\ &= \sum_{j < i} (-1)^i (-1)^j [v_0, \dots, \hat{v}_j, \dots, \hat{v}_i, \dots, v_n] - \sum_{j < i} (-1)^i (-1)^j [v_0, \dots, \hat{v}_j, \dots, \hat{v}_i, \dots, v_n] \\ &= 0 \end{aligned}$$

■

Definition 2.9 Given a simplicial complex X , the n^{th} homology group of X is

$$H_n(X) = \text{Ker } \partial_n / \text{Im } \partial_{n+1},$$

where ∂_n is the boundary homomorphism on $\Delta_{n+1}(X)$.

We give some intuition on a small simplicial complex in Figure 2.5. Intuitively, simplicial homology counts holes in a complex. In this case, there are two holes of dimension 1. Thus, we expect to have two generators for the homology group of dimension 1.

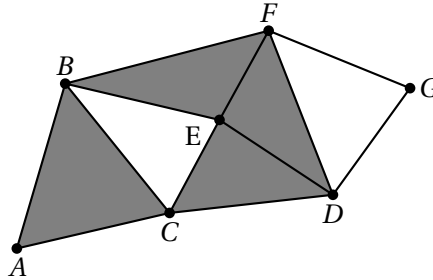


Figure 2.5 – Homology of a small simplicial complex

A representative cycle of the left hole could be the three edges $[BC]$, $[CE]$ and $[EB]$, id est the chain $[BC] + [CE] + [EB]$. For commodity, we write this cycle (BCE) . Remark that the cycle $(ACEB)$ is in the same homology class as (BCE) . The difference between the two of them is $(ACB) = [AC] + [CB] + [BA]$ which is the boundary of a 2-simplex and thus belongs to $\text{Im } \partial_2$. More complex elements are also representatives of this class, $(ACDEFBEFBC)$ for example.

We can also choose a short representative for the right hole by taking (DGF) . Now, the cycle $(BCDGF)$ is in neither of these classes. However, it is equal to $(BEC) + (DGF)$ in the quotient group. These two cycles constitutes a set of generators of the homology group of dimension 1. Counting holes is equivalent to counting generators.

Many properties we prove are true for all homology groups of the simplicial complex X . In this case, we write $H_*(X)$ instead of $H_n(X)$ to denote the fact that we can choose the group arbitrarily. We defined simplicial homology with coefficients in a ring \mathbb{A} . In practice, we restrict ourselves to coefficient in finite fields for computation reasons. From now, \mathbb{A} is assumed to be a finite field.

2.1.4 Simplicial maps

We can relate simplicial complexes using the notion of simplicial maps.

Definition 2.10 *Let X and Y be two simplicial complexes. A simplicial map $f : X \rightarrow Y$ is an application such that for any simplex $\sigma \in X$, $f(\sigma) = \cup_{p \in \sigma} f(p)$ is a simplex of Y . Moreover, two simplicial maps $f : X \rightarrow Y$ and $g : X \rightarrow Y$ are contiguous if $\sigma \in X$ implies that $f(\sigma) \cup g(\sigma) \in Y$.*

Simplicial maps induce homomorphisms between homology groups of X and Y . Moreover, contiguous maps induce the same homomorphisms in simplicial homology. Combining Theorems 12.4 and 12.5 from [83]:

Theorem 2.11 *Two contiguous simplicial maps $f, g : X \rightarrow Y$ induce two homomorphisms f_\star and g_\star that are equal in simplicial homology.*

We use the notion of contiguity for technical results in Chapter 5. We will only consider clique complexes at this time and the contiguity will be proved using the following technical lemma:

Lemma 2.12 *Let X and Y be clique complexes and let f and g be two functions from the vertex set of X to the vertex set of Y . If for every edge $(p, q) \in X$, the simplex generated by $\{f(p), g(p), f(q), g(q)\}$ is in Y , then f and g induce contiguous simplicial maps from X to Y .*

Proof: Let σ be a simplex of X . Every pair in $f(\sigma) \cup g(\sigma)$ is of the form $(f(p), f(q))$, $(f(p), g(q))$, or $(g(p), g(q))$ for some vertices p and q in σ . Since $(p, q) \in \sigma$, the hypothesis of the lemma implies that all of these pairs are edges of Y . Thus, $f(\sigma) \cup g(\sigma)$ is a simplex in Y because Y is a clique complex. Moreover, $f(\sigma) \in Y$ and $g(\sigma) \in Y$ because simplices are closed under taking subsets. Therefore, f and g are contiguous simplicial maps. ■

2.1.5 Singular homology

Objects in real life are not simplicial complexes so analysing them using simplicial homology does not make a lot of sense except if we can relate it to more general objects. This is done by building singular homology. We denote the standard n -simplex by $\Delta^n = \{(t_0, \dots, t_n) \in \mathbb{R}^{n+1} \mid \sum t_i = 1 \wedge t_i \geq 0\}$.

Definition 2.13 *A singular n -simplex in a topological space \mathbb{X} is a continuous map $\sigma : \Delta^n \rightarrow \mathbb{X}$.*

The construction done with abstract and geometric complexes can also be done with singular complexes to define singular homology groups. For simplicial complexes, both constructions are equivalent [69, Theorem 2.27].

Theorem 2.14 *Let X be a simplicial complex. Then the simplicial and singular homology groups of X are isomorphic.*

In this thesis, we restrict ourselves to the class of *triangulable* metric spaces.

Definition 2.15 *A metric space \mathbb{X} is triangulable if it is homeomorphic to a locally finite simplicial complex.*

The homeomorphism between \mathbb{X} and a locally finite simplicial complex C , either abstract or geometric, implies that their singular homology groups are isomorphic. Hence, the singular homology of \mathbb{X} is equivalent to the simplicial homology of C . Therefore, we can compute the homology of C in order to obtain the homology of X . More details on singular homology can be found in [69, 83].

2.2 Persistence

In this section, we define persistence diagrams and study how to compare them. We provide a way to prove the closeness of persistence diagrams using the notion of interleaving between persistence modules. Finally, we quickly discuss the computation of persistence diagrams.

2.2.1 Persistence diagrams

First, we introduce the basic vocabulary of persistent homology. A *filtration* $\{F_\alpha\}_{\alpha \in \mathbb{R}}$ is a family of topological spaces F_α such that for any $\alpha \leq \beta$, $F_\alpha \subset F_\beta$. The inclusion of F_α into F_β is denoted by $F_\alpha \hookrightarrow F_\beta$. A *persistence module* is a family of vector spaces $\{U_\alpha\}_{\alpha \in \mathbb{R}}$ over a field k and of homomorphisms $u_\alpha^\beta : U_\alpha \rightarrow U_\beta$ such that for all $\alpha \leq \beta \leq \gamma$, $u_\alpha^\gamma = u_\beta^\gamma \circ u_\alpha^\beta$ and $u_\alpha^\alpha = Id$. Given a filtration $\mathcal{F} = \{F_\alpha\}_{\alpha \in \mathbb{R}}$ and $\alpha \leq \beta$, the inclusion induces a homomorphism at the homology level $H_*(F_\alpha) \rightarrow H_*(F_\beta)$. These homomorphisms and the homology groups of F_α form one persistence module for each dimension. For the sake of simplicity, we call the set of these modules, the persistence module of \mathcal{F} .

Definition 2.16 *A persistence module $\mathbb{U} = (\{U_\alpha\}, \{u_\alpha^\beta\})$ is quadrant-tame or q -tame if all homomorphisms u_α^β with $\alpha < \beta$ have finite rank. By extension, a filtration $\mathcal{F} = \{F_\alpha\}_{\alpha \in \mathbb{R}}$ is said q -tame if its persistence module is q -tame.*

This notion of tameness ensures some basic properties for the persistence module. For example, a q -tame persistence module has a well-defined persistence diagram [26, Theorem 2.8]. Other kind of tameness exists but we limit ourselves to q -tameness. A complete overview of the tameness properties can be found in [26].

We define the direct sum $\mathbb{W} = \mathbb{U} \oplus \mathbb{V}$ of two persistence modules as $W_\alpha = U_\alpha \oplus V_\alpha$ and $w_\alpha^\beta = u_\alpha^\beta \oplus v_\alpha^\beta$. A persistence module \mathbb{W} is *indecomposable* if the only decompositions of \mathbb{W} are $\mathbb{W} \oplus 0$

and $0 \oplus \mathbb{W}$. Given an interval $J \subset \mathbb{R}$, the interval module \mathbb{I}^J is defined as the persistence module such that $I_t = k$ if $t \in J$ and 0 otherwise and $i_s^t = Id$ if $s, t \in J$ and 0 otherwise.

Interval modules are indecomposable [26, Proposition 1.2] and a persistence module \mathbb{V} where each V_i is finite dimensional can be decomposed into a direct sum of interval modules $\mathbb{V} = \bigoplus_{l \in L} \mathbb{I}^{J_l}$ [42]. Moreover, this decomposition is unique up to a bijection between the intervals [26, Theorem 1.3]. It means that if $\mathbb{V} = \bigoplus_{l \in L} \mathbb{I}^{J_l} = \bigoplus_{m \in M} \mathbb{I}^{K_m}$ then there exists a bijection $\sigma : L \rightarrow M$ such that $J_l = K_{\sigma(l)}$ for all l .

Each interval module \mathbb{I}^J in the decomposition \mathbb{V} can be seen as a homology generator that appears at the start of J and disappears at the end of J . These two points are called respectively *birth* and *death* of the generator. The set of interval modules in the decomposition describes the persistence module and we define the persistence diagram using these intervals.

Definition 2.17 *The persistence diagram $\text{Dgm}(\mathbb{V})$ of a q -tame persistence module $\mathbb{V} = \bigoplus_{l \in L} \mathbb{I}^{J_l}$ is the multi-set:*

$$\text{Dgm}(\mathbb{V}) = \{(p_l, q_l) | l \in L\}$$

where p_l and q_l are the extremities of J_l .

For persistent homology, the number of points at coordinates (a, b) corresponds to the number of generator that appears in the homology groups at time a and disappear at time b . An alternate definition of persistence diagrams rely on counting rectangles. We consider two sequences $A = (a_i)_{i \in \mathbb{Z}}$ and $B = (b_j)_{j \in \mathbb{Z}}$ such that $a_i < a_{i+1}$ and $b_j < b_{j+1}$ for all i and j , $\lim_{i \rightarrow -\infty} a_i = \lim_{j \rightarrow -\infty} b_j = -\infty$ and $\lim_{i \rightarrow \infty} a_i = \lim_{j \rightarrow \infty} b_j = \infty$. Looking at $\text{rk}(H_*(F_{a_i}) \rightarrow H_*(F_{b_j}))$, we have the number of generators that existed in F_{a_i} and still exist in F_{b_j} .

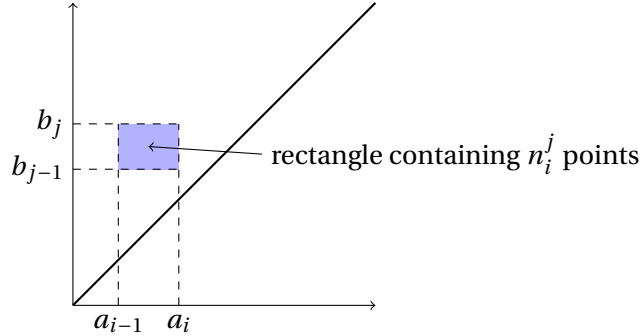


Figure 2.6 – Discrete construction of persistence diagram

Now, using the ranks of the homomorphisms around a_i and b_j , we can obtain the number n_i^j of classes that are born between a_{i-1} and a_i and die between b_{j-1} and b_j for $b_j > a_{i-1}$.

$$\begin{aligned} n_i^j = & \text{rk}(H_k(F_{a_i}) \rightarrow H_k(F_{b_{j-1}})) + \text{rk}(H_k(F_{a_{i-1}}) \rightarrow H_k(F_{b_j})) \\ & - \text{rk}(H_k(F_{a_i}) \rightarrow H_k(F_{b_j})) - \text{rk}(H_k(F_{a_{i-1}}) \rightarrow H_k(F_{b_{j-1}})) \end{aligned}$$

This number n_i^j corresponds to the number of points present in the blue rectangle of Figure 2.6. The counting rectangles can be refined by building two sequences of sequences $A_n = (a_i^n)_{i \in \mathbb{Z}}$ and $B_n = (b_j^n)_{j \in \mathbb{Z}}$ such that $\limsup_{i,j} \{|b_j - b_{j-1}|; |a_i - a_{i-1}|\} = 0$. Both constructions of persistence diagrams are equivalent for q -tame persistence modules and points outside the diagonal $\Delta = \{(x, x) | x \in \mathbb{R}\}$. Points located on Δ are not useful for inference purposes. They are ignored by the notion of distance used to compare persistence diagrams.

A persistence diagram can be represented in different ways. We use two equivalent representations in this thesis. Figure 2.7 shows these two representations for a given diagram. The first one is the most direct. Seeing the diagram as a multi-point set of \mathbb{R}^2 , we just draw each point on the plane. While this representation is natural and is nice to intuitively comprehend the notion of distance that will be introduced later, it does not work well when there exist points with multiplicity. The second way of drawing the diagram is called a *barcode*. Each point is now represented by an interval, which starts at its birth and stops at its death. Sometimes less compact, this representation works better when points have multiplicity more than one as we draw as many intervals as their multiplicity. In this example, two of the points have multiplicity 2, fact that is not visible on the first representation.

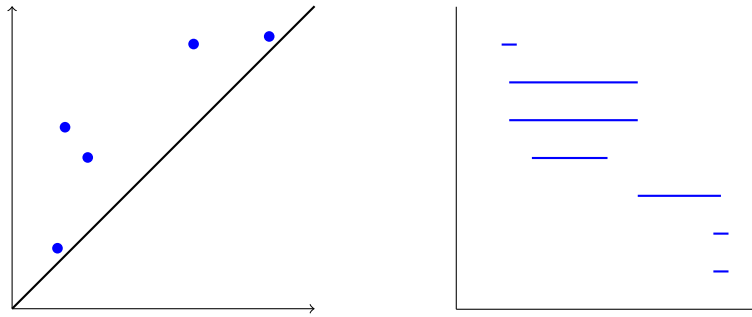


Figure 2.7 – Two representations of the same diagram

In this thesis, we often build filtrations using functions. By abuse of notation, we speak of the filtration of a function instead of its sub-level sets filtration $\{f^{-1}([-\infty, \alpha])\}_{\alpha \in \mathbb{R}}$. By extension, we speak about the persistence diagram of f and consider the tameness of f instead of the ones of its sub-level sets filtration.

2.2.2 Bottleneck distance

To compare persistence diagram, we use the bottleneck distance. The idea is that two persistence diagrams are close if their features with long lifespan have close birth and death times. First we introduce the notion of δ -matching.

Definition 2.18 *Let P and Q be two multi-sets of points in \mathbb{R}^2 . A δ -matching between P and Q is a collection of pairs $M \subset P \times Q$ such:*

- *Any point of P is matched with at most one point of Q and reciprocally.*
- $\forall (p, q) \in M, \|p - q\|_\infty \leq \delta.$

- If $a \in P \cup Q$ is unmatched, then a is at distance at most δ from the diagonal Δ for $\|\cdot\|_\infty$.

Definition 2.19 Let P and Q be two multi-sets of points in \mathbb{R}^2 . The bottleneck distance between P and Q is defined by:

$$d_B(P, Q) = \inf\{\delta \mid \exists M \text{ a } \delta\text{-matching between } P \text{ and } Q\}$$

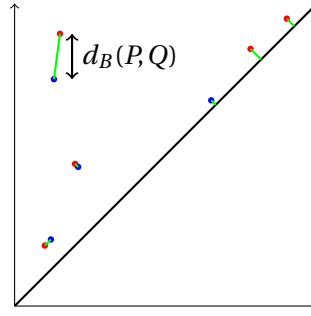


Figure 2.8 – Bottleneck distance

The bottleneck distance ignores points that are close to the diagonal. In the persistent homology setting, it means that only the elements with a long lifespan have to be matched for two diagrams to be close.

We also use a variant of the bottleneck distance. When comparing sets which do not have the same scale, we use the bottleneck in a logarithmic scale. Consider an object K sampled with n points P . We build the point cloud P' using a homothety of factor 10. Then P' sampled an object K' which is ten times the size of K . The bottleneck distance between the diagrams of P and P' is large. However, the topologies of K and K' are the same but at different scales. By looking at the points in a logarithmic scale, the persistence diagrams are much more similar as shown in Figure 2.9. On both diagrams, blue points correspond to the diagram of the small object and red squares to the one of the big object. In logarithmic scale, they are identical up to a translation along the first bisector, which corresponds to the scale factor.

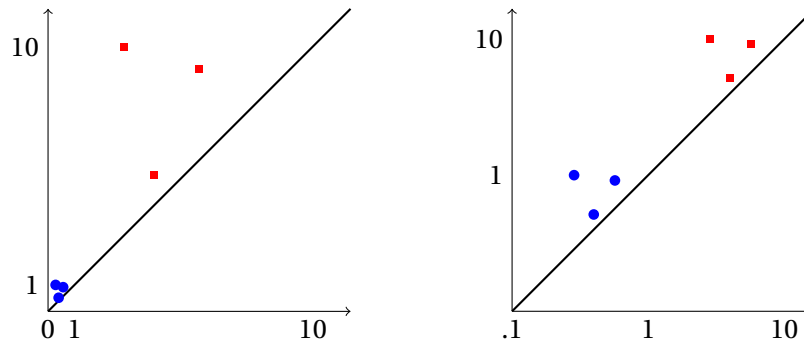


Figure 2.9 – Comparison of persistence diagrams of an object at two different scales

Definition 2.20 Let P and Q be two multi-sets of points in \mathbb{R}^2 . Consider the multi sets P^l and Q^l images of P and Q by the application $(x, y) \mapsto (\ln(x), \ln(y))$. The bottleneck distance in logarithmic scale between P and Q is defined by:

$$d_B^{\log}(P, Q) = d_B(P^l, Q^l)$$

2.2.3 Persistence module interleaving

In practice, we compare persistence modules using the notion of interleaving. Theorem 2.23 then implies a proximity of the persistence diagrams for the bottleneck distance. First, we introduce the notion of ϵ -homomorphism.

Definition 2.21 Let $\mathbb{U} = (\{U_\alpha\}, \{u_\alpha^\beta\})$ and $\mathbb{V} = (\{V_\alpha\}, \{v_\alpha^\beta\})$ be two persistence modules. An ϵ -homomorphism from \mathbb{U} to \mathbb{V} is a collection $\Phi = \{\phi_t : U_t \rightarrow V_{(t+\epsilon)}\}_{t \in \mathbb{R}}$ of linear maps such that the diagram of Figure 2.10 commutes for all $\alpha \leq \beta$.

$$\begin{array}{ccc} U_\alpha & \xrightarrow{u_\alpha^\beta} & U_\beta \\ \phi_\alpha \downarrow & & \downarrow \phi_\beta \\ V_{\alpha+\epsilon} & \xrightarrow{v_{\alpha+\epsilon}^{\beta+\epsilon}} & V_{\beta+\epsilon} \end{array}$$

Figure 2.10 – Diagram of an ϵ -homomorphism

ϵ -homomorphisms are sometime called homomorphisms of degree ϵ . The collection of maps $\{u_t^{t+\epsilon}\}$ is an ϵ -endomorphism and is called the shift map $1_{\mathbb{U}}^\epsilon$. The notion of ϵ -interleaving relates two persistence modules by guaranteeing that we can go from one to the other and then come back. It is not an inverse as it does not give us the identity, but the shift map $1^{2\epsilon}$.

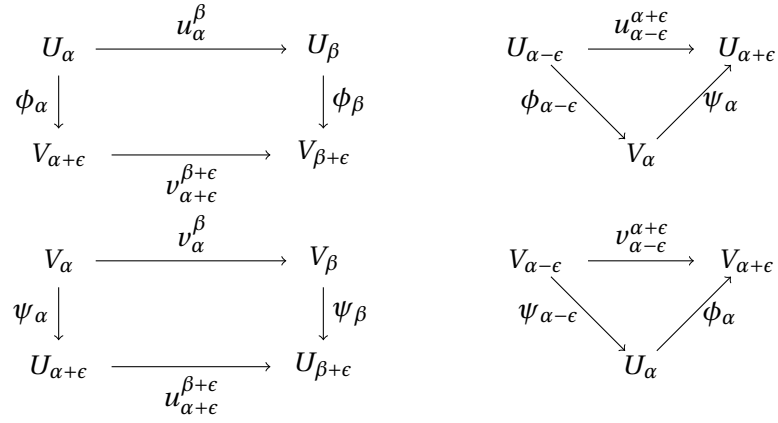
Definition 2.22 Let $\mathbb{U} = (\{U_\alpha\}, \{u_\alpha^\beta\})$ and $\mathbb{V} = (\{V_\alpha\}, \{v_\alpha^\beta\})$ be two persistence modules. \mathbb{U} and \mathbb{V} are ϵ -additively interleaved if there exists two ϵ -homomorphism Φ from \mathbb{U} to \mathbb{V} and Ψ from \mathbb{V} to \mathbb{U} such that:

$$\Psi\Phi = 1_{\mathbb{V}}^{2\epsilon} \quad \Phi\Psi = 1_{\mathbb{U}}^{2\epsilon}$$

The definition of ϵ -interleaving is equivalent to the commutativity of the diagrams in Figure 2.11. This notion of interleaving is called additive due to the addition of parameter ϵ when doing shifts in the diagram and by opposition to the multiplicative interleaving defined later. Unless stated otherwise, we speak of two ϵ -interleaved modules when the interleaving is additive.

The main stability result for persistence diagrams states that two interleaved q -tame persistence modules have diagrams close with respect to the bottleneck distance [22, 26].

Theorem 2.23 Let $\mathbb{U} = (\{U_\alpha\}, \{u_\alpha^\beta\})$ and $\mathbb{V} = (\{V_\alpha\}, \{v_\alpha^\beta\})$ be two q -tame ϵ -interleaved persistence modules, then $d_B(\text{Dgm}(\mathbb{U}), \text{Dgm}(\mathbb{V})) \leq \epsilon$.

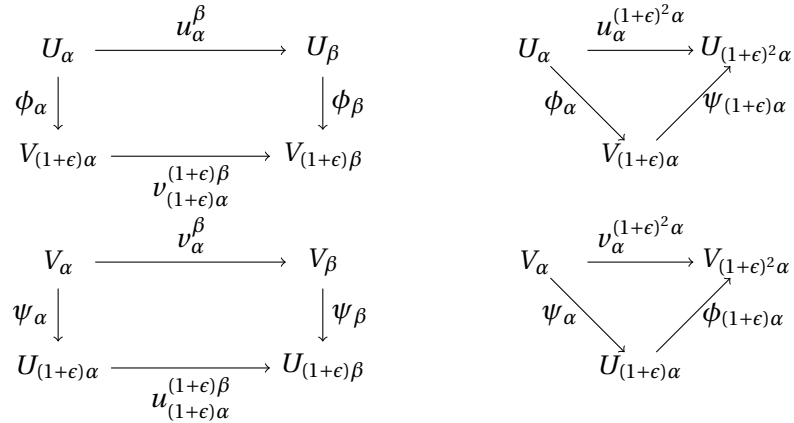

 Figure 2.11 – Commutative diagrams of additive ϵ -interleaving

The proof of this theorem is quite long and technical. The interested reader can either refer to the complete proofs in [22, 26] or the one restricted to sub-level sets filtrations of q -tame functions in [40].

It is not always possible to achieve an ϵ -interleaving in practice but we can sometimes work in a logarithmic scale. Given a persistence module $\mathbb{U} = (\{U_\alpha\}, \{u_\alpha^\beta\})$, we consider the persistence module $\mathbb{U}' = (\mathbb{U}_{\ln(\alpha)}, \{u_{\ln(\alpha)}^{\ln(\beta)}\})$.

Definition 2.24 Let $\mathbb{U} = (\{U_\alpha\}, \{u_\alpha^\beta\})$ and $\mathbb{V} = (\{V_\alpha\}, \{v_\alpha^\beta\})$ be two persistence modules. Let \mathbb{U}' and \mathbb{V}' be the persistence modules considered in logarithmic scale. \mathbb{U} and \mathbb{V} are ϵ -multiplicatively interleaved if \mathbb{U}' and \mathbb{V}' are ϵ -additively interleaved.

As for the additive interleaving, it corresponds to a set of commutative diagrams given in Figure 2.12


 Figure 2.12 – Commutative diagrams of multiplicative ϵ -interleaving

The use of logarithmic scale immediately implies:

Corollary 2.25 *Let $\mathbb{U} = (\{U_\alpha\}, \{u_\alpha^\beta\})$ and $\mathbb{V} = (\{V_\alpha\}, \{v_\alpha^\beta\})$ be two multiplicatively ϵ -interleaved persistence modules, then:*

$$d_B^{\log}(\text{Dgm}(\mathbb{U}), \text{Dgm}(\mathbb{V})) \leq \ln(1 + \epsilon).$$

Proof: The persistence modules \mathbb{U} and \mathbb{V} can be looked at logarithmic scale to obtain \mathbb{U}' and \mathbb{V}' that are additively $\ln(1 + \epsilon)$ -interleaved. By definition of the logarithmic bottleneck distance, $d_B^{\log}(\text{Dgm}(\mathbb{U}), \text{Dgm}(\mathbb{V})) = d_B(\text{Dgm}(\mathbb{U}'), \text{Dgm}(\mathbb{V}'))$. Using Theorem 2.23, we obtain

$$d_B^{\log}(\text{Dgm}(\mathbb{U}), \text{Dgm}(\mathbb{V})) = d_B(\text{Dgm}(\mathbb{U}'), \text{Dgm}(\mathbb{V}')) \leq \ln(1 + \epsilon)$$

■

Sometimes, interleaving is done directly on filtrations instead of persistence modules.

Definition 2.26 *Let $\{U_\alpha\}_{\alpha \in \mathbb{R}}$ and $\{V_\alpha\}_{\alpha \in \mathbb{R}}$ be two filtrations and $\epsilon \geq 0$. $\{U_\alpha\}$ and $\{V_\alpha\}$ are ϵ -interleaved if for any $\alpha \in \mathbb{R}$, $U_\alpha \subset V_{\alpha+\epsilon} \subset U_{\alpha+2\epsilon}$.*

Remark that it directly implies that their persistence modules are interleaved.

Corollary 2.27 *Let $\{U_\alpha\}_{\alpha \in \mathbb{R}}$ and $\{V_\alpha\}_{\alpha \in \mathbb{R}}$ be two q -tame and ϵ -interleaved filtrations. Then $d_B(\text{Dgm}(\{U_\alpha\}), \text{Dgm}(\{V_\alpha\})) \leq \epsilon$.*

Proof: The two filtrations $\{U_\alpha\}$ and $\{V_\alpha\}$ with the inclusions induce two persistence modules \mathbb{U} and \mathbb{V} . $\{U_\alpha\}$ and $\{V_\alpha\}$ are ϵ -interleaved. Hence, the diagram 2.11 commutes when all maps are inclusions, which means that \mathbb{U} and \mathbb{V} are ϵ -interleaved and Theorem 2.23 applies. ■

2.2.4 Computation of persistence diagrams

A *filtered simplicial complex* is a simplicial complex S with a function $f : S \rightarrow \mathbb{R}$ such that for any pair of simplices σ, σ' of S such that $\sigma' \subset \sigma$, $f(\sigma') \leq f(\sigma)$. We have a filtration $\mathcal{F} = \{F_\alpha\}$ where $F_\alpha = f^{-1}([-\infty, \alpha])$ is a simplicial complex for any $\alpha \in \mathbb{R}$. This filtration induces a persistence module. Taking the homology coefficients in a finite field, the computation of $\text{Dgm}(\mathcal{F})$ is equivalent to the reduction of a matrix of size $N \times N$, where N is the number of simplices in F_∞ [57]. The worst case complexity is thus $O(N^3)$.

However, this matrix is usually sparse. Algorithms have been developed to improve the running time of persistence diagram computation [8, 14, 99]. Versions using the inherent duality between homology and cohomology also exist [12, 44, 45]. In practice, these algorithms have a running time that is linear in the number of simplices N .

Until now, we restricted ourselves to the computation of persistent homology for filtrations of simplicial complexes. This will be sufficient for the work in this thesis. It should however be noted that the definition of persistence modules is not restricted to sequences of sets that are filtrations. Transposed to simplicial complexes, this means that we can discuss sequences where we have simplicial maps between complexes. It was first introduced for the *zig-zag persistence* [19] where inclusions can occur in both directions, i.e. the family $\{F_i\}_{i \in \mathbb{N}}$ verifies for every i , either $F_i \subset F_{i+1}$ or $F_i \supset F_{i+1}$. Algorithms to compute this kind of diagrams have been developed [20, 81] and then extended to any simplicial maps [48].

2.3 Geometry

Algebraic preliminaries were necessary to introduce the concept of persistence diagrams. However, guarantees on the result of algorithms for topological data analysis rely on geometric assumptions. In this section, we make a quick overview of some common geometric concepts and results that will be used during the thesis.

2.3.1 Metric spaces

We mostly work with metric spaces. Recall that a metric space \mathbb{X} is a set with a distance function $d_{\mathbb{X}}(\cdot, \cdot)$ defined for every pair of elements. The distance is a non-negative and symmetric function that satisfies the triangle inequality, i.e. $d_{\mathbb{X}}(x, y) \leq d_{\mathbb{X}}(x, z) + d_{\mathbb{X}}(y, z)$ for any x, y, z and such that $d_{\mathbb{X}}(x, y) = 0$ if and only if $x = y$. Given an element $x \in \mathbb{X}$ and a real radius $r \geq 0$, we denote $B(x, r) = \{y \in \mathbb{X} \mid d_{\mathbb{X}}(x, y) < r\}$ the open ball of centre x and radius r . Similarly, the closed ball will be denoted $\bar{B}(x, r) = \{y \in \mathbb{X} \mid d_{\mathbb{X}}(x, y) \leq r\}$.

Subsets of \mathbb{X} can be compared against each other using the Hausdorff distance. Given P and Q two subsets of \mathbb{X} , the idea is to measure how far a point of P has to be moved before being in Q , and reciprocally. Formally, the Hausdorff distance is given by:

$$d_H(P, Q) = \max \left(\sup_{y \in Q} \inf_{x \in P} d_{\mathbb{X}}(x, y); \sup_{x \in P} \inf_{y \in Q} d_{\mathbb{X}}(x, y) \right).$$

On the space of compact subsets of \mathbb{X} , d_H verifies the axioms of a distance. Introducing the function distance to the subset P as $d_P(x) = \inf_{p \in P} d_{\mathbb{X}}(x, p)$ and the distance d_Q to Q , the Hausdorff distance can be rewritten as the infinite norm between d_P and d_Q . Furthermore, it can be characterised using inclusions of sub-level sets.

Lemma 2.28 *Given two subsets P and Q of a metric space \mathbb{X} , $d_H(P, Q) = \|d_P - d_Q\|_{\infty}$.*

Proof: Let x be a point of \mathbb{X} . For any $\epsilon > 0$, there exists a point $y \in P$ such that $d_{\mathbb{X}}(x, y) \leq d_P(x) + \epsilon$. Moreover, there exists $z \in Q$ such that $d_{\mathbb{X}}(y, z) \leq d_H(P, Q) + \epsilon$. Hence $d_Q(x) \leq d_P(x) + d_H(P, Q) + 2\epsilon$. This is true for all $\epsilon > 0$ and therefore $\|d_Q - d_P\|_{\infty} \leq d_H(P, Q)$.

There also exists a point $x' \in \mathbb{X}$ such that $|d_Q(x') - d_P(x')| \geq \|d_Q - d_P\|_{\infty} - \epsilon$. Without loss of generality, we assume that $d_P(x') \leq d_Q(x')$. There exists $y' \in P$ such that $d_{\mathbb{X}}(x', y') \leq d_P(x') + \epsilon$. For any $z' \in Q$, $d_{\mathbb{X}}(y', z') \geq d_{\mathbb{X}}(x', z') - d_{\mathbb{X}}(x', y') \geq d_Q(x') - d_P(x') - \epsilon$. Hence $d_H(P, Q) \geq \|d_Q - d_P\|_{\infty}$. ■

Lemma 2.29 *Given two subsets P and Q of a metric space \mathbb{X} , the Hausdorff distance $d_H(P, Q)$ is less than ϵ if and only if, for any $\alpha \in \mathbb{R}$,*

$$d_P^{-1}([-\infty, \alpha]) \subset d_Q^{-1}([-\infty, \alpha + \epsilon]) \subset d_P^{-1}([-\infty, \alpha + 2\epsilon]).$$

Proof: Remark that if $\alpha < 0$, then $d_P^{-1}([-\infty, \alpha]) = \emptyset$. Assume that $d_H(P, Q) \leq \epsilon$, take $\alpha \geq 0$ and let x be a point of $d_P^{-1}([-\infty, \alpha])$. Then for any $\eta > 0$, there exists a point $p \in P$ such that $d_{\mathbb{X}}(x, p) \leq \alpha + \eta$. By definition of the Hausdorff distance, there exists a point $q \in Q$ such that $d_{\mathbb{X}}(p, q) \leq \epsilon + \eta$. Using the triangle inequality, $d_{\mathbb{X}}(x, q) \leq \alpha + \epsilon + 2\eta$. As there exists a $q \in Q$ for

any value of η , this means that $d_Q(x) \leq \alpha + \epsilon$. The second inequality is obtained by reversing the roles of P and Q .

Now assume that $d_P^{-1}([-\infty, \alpha]) \subset d_Q^{-1}([-\infty, \alpha + \epsilon]) \subset d_P^{-1}([-\infty, \alpha + 2\epsilon])$. This implies that $P \subset d_P^{-1}([-\infty, 0]) \subset d_Q^{-1}([-\infty, \epsilon])$. In other words, for any $p \in P$ and $\eta > 0$, there exists a $q \in Q$ such that $d_{\mathbb{X}}(p, q) \leq \epsilon + \eta$. Hence, $\sup_{p \in P} \inf_{q \in Q} d_{\mathbb{X}}(p, q) \leq \epsilon$. Doing the same with the second inequality and $\alpha = -\epsilon$, we obtain $d_H(P, Q) \leq \epsilon$. ■

Subsets of a metric space can have an intrinsic dimension that is smaller than the dimension of the ambient space. For example, a curve in the Euclidean plane is of intrinsic dimension 1, while the ambient space is of dimension 2. The intrinsic dimension is described by the doubling dimension.

Definition 2.30 *The doubling constant $\lambda_{\mathbb{X}}$ of a metric space \mathbb{X} is the maximum over all balls $B(x, r)$ with $x \in \mathbb{X}$ of the minimum number of balls of radius $r/2$ required to cover $B(x, r)$. The doubling dimension is defined to be $\log_2(\lambda_{\mathbb{X}})$.*

When working in Euclidean spaces, we sometimes encounter the assumption that an input point set is in general position.

Definition 2.31 *Let P be a point set in an Euclidean space \mathbb{R}^d . P is in general position if for any $i < d$, no set of $i + 3$ distinct points of P is on a sphere of dimension i and the convex hull of any $d + 1$ distinct points of P is of dimension d .*

This condition means for example that no triple of points are aligned and no quadruple of points are co-circular. Some geometric constructions such as the Delaunay triangulation requires this condition to avoid degenerate cases that can hinder the complexity or the soundness of algorithms. The general position assumption is reasonable if the point set P is finite. If P is not in general position, it suffices to slightly perturb every point and we obtain a point set that is in general position with probability 1.

2.3.2 Compact sets

Usually, we are given an input point set P and we assume the existence of an underlying object K . Our objective is to study K , sometimes called the *ground truth* using P . Assumptions on the Hausdorff distance between P and K are common to provide theoretical guarantees. The distance usually needs to be bounded by some geometric quantities describing K . Here, we consider the case of a compact set K .

Every $x \in \mathbb{X}$ has a closest point in K . However, this point is not necessarily unique. The set of all points that have at least two nearest neighbours on K is called the *medial axis* of K and is illustrated in Figure 2.13.

Definition 2.32 *Let K be compact set of metric space \mathbb{X} and A its medial axis. The reach of K is defined as*

$$r_K = \inf_{x \in K} d_{\mathbb{X}}(x, A)$$

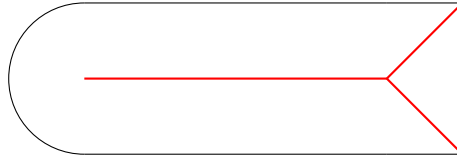


Figure 2.13 – A curve and its medial axis

A standard assumption for reconstruction is that the reach of K is not too small. For example, it has to be positive or more than the Hausdorff distance between K and the point set P . Intuitively, this ensures that we have no point close to the medial axis and hence, given a point $p \in P$, we can approximate well where its projection P is located.

The use of the reach has some limitations when there exist some sharp vertices for example. The reach of the surface of a cube is 0. More precise quantities have been proposed for sampling compact sets such as the weak feature size and the μ -reach [24], however we only use the notion of reach in this thesis.

2.3.3 Riemannian geometry

Riemannian manifolds have nice regularity properties and are often used to state hypotheses. Here, we introduce smooth Riemannian submanifolds of Euclidean spaces and some concepts we need to express our assumptions. For a more thorough presentation, the reader can refer to specialised books such as [52, 63].

Definition 2.33 A smooth d' -submanifold of \mathbb{R}^d is a subset $M \subset \mathbb{R}^d$ such that for any point $x \in M$, there exists an open set U containing x and a C^∞ diffeomorphism ϕ from U onto an open subset $\phi(U) \subset \mathbb{R}^d$ such that $\phi(U \cap M)$ is a vector subspace of dimension d' .

The scalar product on the tangent spaces of a submanifold M of \mathbb{R}^d induces a metric d_M , called the *geodesic metric*. The pair (M, d_M) is a metric space called a *Riemannian submanifold* of \mathbb{R}^d . As for any metric space, we can speak of balls defined on this space. The shortest path between two points x and y is called the *minimizing geodesic*. Remark that the geodesic distance between two points on a submanifold of \mathbb{R}^d is always greater than their distance in the metric of the embedding space \mathbb{R}^d . The shortest path to go from one point to another is always longer when we restrict ourselves to stay on a subspace rather than being allowed to move in the whole space.

The data we can work with is given under the form of point sets. We usually assume that the point set samples the manifold M with a certain precision.

Definition 2.34 Let M be a Riemannian manifold and let $P \subset M$ be a point set. P is a Riemannian ϵ -sampling of M if for any $x \in M$, there exists $p \in P$ such that $d_M(x, p) \leq \epsilon$.

Remark that if P is a Riemannian ϵ -sampling of M , then $d_H(P, M) \leq \epsilon$. One part of the relation is given by the inclusion of P inside M and the other is directly implied by the sampling assumption and the fact that the Riemannian metric is greater than the Euclidean metric.

Riemannian manifolds are very general objects and we cannot hope to reconstruct them or infer correct information on them without adding some regularity assumptions.

To describe the way a Riemannian submanifold M folds, we use the notion of *sectional curvature* at a point p . It is a function from the space of tangent 2-planes of M at p to \mathbb{R} , whose value is the curvature of a geodesic starting at p and tangent to the plane. This notion of curvature is an extension to higher dimensions of the usual curvature notion for curves and surfaces. For a formal definition, we refer the reader to [52, Chapter 4]. The sectional curvature being the only notion of curvature used in this thesis, it will simply be called curvature.

Proposition 2.35 ([49, Proposition 2.1]) *Let $M \subset \mathbb{R}^2$ be a smooth compact manifold without boundary. Then the absolute value of the curvature of M is bounded by c_M which verifies:*

$$c_M \leq \frac{2}{r_M^2}$$

where r_M is the reach of M .

Another useful concept is the *strong convexity radius* $\rho(M)$ defined as the largest radius such that for any point of $x \in M$ and radius $r < \rho(M)$, the geodesic ball $\mathcal{B}(x, r)$ is strongly convex, i.e. the shortest geodesic between two points y and z of $\mathcal{B}(x, r)$ is unique and enclosed inside $\mathcal{B}(x, r)$.

2.3.4 Lower bound on the volume of Riemannian balls

Assuming we sample a compact manifold uniformly, we want to guarantee that we obtain a dense enough point set on the manifold. We obtain it using a lower bound on the volume of Riemannian balls, id est the balls in the Riemannian metric. Consider a manifold M with sectional curvature upper bounded by c_M . Then for any point $x \in M$, the Günther-Bishop theorem provides a lower bound of the volume of the Riemannian ball of radius a .

Theorem 2.36 (Günther-Bishop) *Assuming that the sectional curvature of a manifold M is always less than c_M and a is less than the strong convexity radius of M , then for any point $x \in M$, the volume $\mathcal{V}(x, a)$ of the geodesic ball centred on x and of radius a is greater than $V_{d'}^{c_M}(a)$ where d' is the intrinsic dimension of M and $V_{d'}^{c_M}(a)$ is the volume of the Riemannian ball of radius a on a surface with constant curvature c_M .*

We explicitly bound the value of $\mathcal{V}(x, a)$, with the following technical lemma:

Lemma 2.37 *Let M be a Riemannian manifold with curvature upper bounded by c_M , then for any $x \in M$ and $a \leq \min(\rho(M); \frac{\pi}{\sqrt{c_M}})$, the volume $\mathcal{V}(x, a)$ of the geodesic ball centred at x and of radius a verifies:*

$$\mathcal{V}(x, a) \geq \mathcal{C}_{d'}^{c_M} a^{d'}$$

where $\mathcal{C}_{d'}^{c_M}$ is a constant independent of x and a .

Proof: Given $a \leq \min(\rho(M), \frac{\pi}{\sqrt{c_M}})$, we want to bound the volume $V_{d'}^{c_M}(a)$. Consider the sphere of dimension d' and curvature c_M . The surface $S_{c_M}^{d'-1}$ of the border of a ball of radius $a \leq \frac{\pi}{\sqrt{c_M}}$ on this sphere is given by [66]:

$$S_{c_M}^{d'-1}(a) = 2\Gamma\left(\frac{1}{2}\right)^{d'} \Gamma\left(\frac{d'}{2}\right)^{-1} c_M^{-\frac{1}{2}(d'-1)} \sin^{d'-1}(c_M a)$$

We can bound the value of $V_{d'}^{c_M}(a)$:

$$\begin{aligned} V_{d'}^{c_M}(a) &= \int_0^a S^{d'-1}(l) dl \\ &= \int_0^a 2\Gamma\left(\frac{1}{2}\right)^{d'} \Gamma\left(\frac{d'}{2}\right)^{-1} c_M^{-\frac{1}{2}(d'-1)} \sin^{d'-1}(c_M l) dl \\ &\geq 2\Gamma\left(\frac{1}{2}\right)^{d'} \Gamma\left(\frac{d'}{2}\right)^{-1} c_M^{-\frac{1}{2}(d'-1)} 2 \int_0^{\frac{a}{2}} \left(\frac{2c_M l}{\pi}\right)^{d'-1} dl \\ &= 4\Gamma\left(\frac{1}{2}\right)^{d'} \Gamma\left(\frac{d'}{2}\right)^{-1} c_M^{-\frac{1}{2}(d'-1)} \frac{\pi}{2c_M} \int_0^{\frac{c_M a}{\pi}} u^{d'-1} du \end{aligned}$$

Writing

$$\mathcal{C}_{d'}^{c_M} = \frac{4}{d'} \Gamma\left(\frac{1}{2}\right)^{d'} \Gamma\left(\frac{d'}{2}\right)^{-1} \left(\frac{\sqrt{c_M}}{\pi}\right)^{d'-1},$$

and using the Günther-Bishop Theorem, we have for any $a \leq \min(\rho(M), \frac{\pi}{\sqrt{c_M}})$ and any $x \in M$,

$$\mathcal{V}(x, a) \geq \mathcal{C}_{d'}^{c_M} a^{d'}.$$

■

3 Distance to a measure

The limitations of the distance to a compact and the distance to a point set led to the introduction of the *distance to a measure*. First proposed by Chazal, Cohen-Steiner and Mérigot in [25], this function aimed at adding robustness with regard to aberrant values, also called *outliers*. In this chapter, we introduce the function and extend its stability results to any metric space along with some easy but interesting properties. Then we study the identifiability of measures knowing their distances and provide new qualitative results.

3.1 Measures

The original idea behind distance to a measure functions was to consider point sets as measures when they were previously considered as compact sets. Measures are functions defined over a specific algebraic construction called σ -algebra. Intuitively, a measure tries to capture the size of an object. To be coherent with this idea, the measure needs some basic properties.

Definition 3.1 A σ -algebra \mathcal{A} on a set X is a non-empty collection of sets that is stable by complement and countable union.

When we measure objects thanks to a function μ , we want μ to have some natural properties. Especially, if objects are disjoint, the size of their union has to be equal to the sum of their individual sizes. This is called the σ -additivity of μ and is formally defined for any countable family $(U_n)_{n \in \mathbb{N}}$ family of pairwise disjoint sets of \mathcal{A} by $\mu(\bigcup_n U_n) = \sum_n \mu(U_n)$.

Definition 3.2 A measure μ over a σ -algebra \mathcal{A} is a non-negative and σ -additive function defined on \mathcal{A} such that $\mu(\emptyset) = 0$

In this thesis, we will always work in a metric space \mathbb{X} and consider the σ -algebra generated by the open balls of \mathbb{X} . It means that the σ -algebra \mathcal{A} will be the smallest set $\mathfrak{B}(\mathbb{X})$ that contains all open balls and is closed by complement and numerable union. The elements of $\mathfrak{B}(\mathbb{X})$ are called Borel subsets of \mathbb{X} . All measures in this thesis will be defined over $\mathfrak{B}(\mathbb{X})$ and, by abuse of notation, we will write that measures are defined on the metric space \mathbb{X} . More specific kinds of measures interest us. First of all is the notion of *probability measure*.

Definition 3.3 A probability measure on a metric space \mathbb{X} is a measure such that $\mu(\mathbb{X}) = 1$.

Given a Borel subset $A \in \mathfrak{B}(\mathbb{X})$, the value $\mu(A)$ is sometimes called the *mass* of A and $\mu(\mathbb{X})$ is called the *total mass* of μ . Thus, a probability measure μ is a measure of total mass 1. Introducing the *empirical measure*, one can consider a set of points as a measure.

Definition 3.4 Let P be a finite set of points in a metric space \mathbb{X} . The empirical measure μ_P is defined as:

$$\mu_P = \frac{1}{|P|} \sum_{p \in P} \delta_p$$

where δ_p is the Dirac measure at point p .

Recall that the Dirac measure at a point p is the measure such that for any set A , $\delta_p(A)$ is equal to 1 if $p \in A$ and 0 otherwise. The empirical measure is the probability measure sharing its mass evenly between all the points of the point set. Remark that in the case of a multi-set, when points can appear more than once, the definition still stands and the points have a mass proportional to their multiplicity. Empirical measures belongs to a larger class of measures called the *measures with finite support*.

Definition 3.5 Given a measure μ defined on a metric space \mathbb{X} , the support of μ , $\text{Supp}(\mu)$, is the smallest closed set of \mathbb{X} whose complement has measure 0.

Definition 3.6 A measure μ is said to have finite support if $\text{Supp}(\mu)$ is a finite point set P .

Measures with finite support and especially empirical measures are often used to approximate other measures. Some discussion can be found in Section 6.3.3.

Definition 3.7 Given a probability measure μ defined on a metric space \mathbb{X} and $1 \leq p \leq \infty$, the p^{th} moment of μ for a point $x_0 \in \mathbb{X}$ is defined by

$$\int_{\mathbb{X}} d_{\mathbb{X}}(x, x_0)^p d\mu(x).$$

3.2 Wasserstein distances

To compare measures with same total mass, we use Wasserstein distances, also called optimal transportation distances. Intuitively, they are the minimal cost to move all the mass from one measure to another. The way two move the mass between two measures is described by a *transport plan*:

Definition 3.8 Let μ and ν be positive measures with the same total mass on a metric space \mathbb{X} . A transport plan between μ and ν is a measure π on $\mathbb{X} \times \mathbb{X}$ such that for all $A, B \in \mathfrak{B}(\mathbb{X})$,

$$\pi(A \times \mathbb{X}) = \mu(A) \text{ and } \pi(\mathbb{X} \times B) = \nu(B).$$

We denote by $\Pi(\mu, \nu)$ the set of all transport plans between μ and ν . The p^{th} order cost of the transport plan π is defined as

$$C_p(\pi) = \left(\int_{\mathbb{X} \times \mathbb{X}} d_{\mathbb{X}}(x, y)^p d\pi(x, y) \right)^{\frac{1}{p}}.$$

The Wasserstein distance between μ and ν is the minimum cost over all transport plans.

Definition 3.9 Let μ and ν be positive measures with the same total mass on a metric space \mathbb{X} . The Wasserstein distance of order p between μ and ν is defined as

$$W_p(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \left(\int_{\mathbb{X} \times \mathbb{X}} d_{\mathbb{X}}(x, y)^p d\pi(x, y) \right)^{\frac{1}{p}}.$$

The Wasserstein distance is finite if both probability measures have finite p -moments, which is always the case for measures with compact support. In this thesis, we will consider the distance W_2 in cases where it is finite. For more details on Wasserstein distances and optimal transportation theory, we refer the reader to [94].

3.3 Distance to a measure

The distance to a measure introduced in [25] is defined in Euclidean spaces. However, the definition is more general and can be extended to any metric space. First, we consider the so-called *pseudo-distance* to a measure.

Definition 3.10 Let μ be a probability measure on a metric space \mathbb{X} and let $m \in [0, 1[$ be a mass parameter. We define the pseudo-distance to μ :

$$\delta_{\mu, m} : x \in \mathbb{X} \mapsto \inf\{r \geq 0 \mid \mu(\bar{B}(x, r)) > m\}.$$

This function corresponds to the distance we need to look at in order to catch the mass m as shown in Figure 3.1. When $m = 0$, it is the distance to the support of μ . For $m > 0$, we not only consider the nearest point of the support but also the structure of the measure beyond.

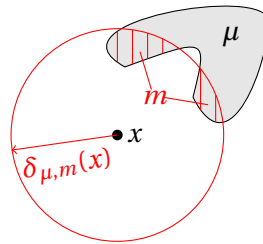


Figure 3.1 – Pseudo-distance to a measure

However, this function lacks essential properties such as the continuity with regard to μ and m . Consider for example the measure μ_ϵ defined on \mathbb{R} and given in Figure 3.2, where 0 has mass $\frac{1}{2} - \epsilon$ and 1 has mass $\frac{1}{2} + \epsilon$. Then $\delta_{\mu_\epsilon, m}(0) = 0$ for any $m < \frac{1}{2} - \epsilon$ and $\delta_{\mu_\epsilon, m}(0) = 1$ for $m \geq \frac{1}{2} - \epsilon$. Moreover, for any $\epsilon > 0$, $\delta_{\mu_\epsilon, \frac{1}{2}}(1) = 0$ but $\delta_{\mu_0, \frac{1}{2}}(1) = 1$, while $W_2(\mu_\epsilon, \mu_0) = \epsilon$. The distance to the measure μ averages the function $m \mapsto \delta_{\mu, m}$ to provide these properties.

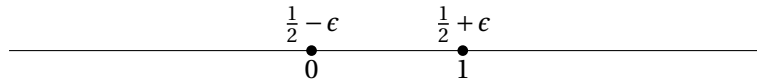


Figure 3.2 – Measure μ_ϵ

Definition 3.11 Let μ be a probability measure on a metric space \mathbb{X} and let $m \in [0, 1[$ be a mass parameter. The distance to the measure μ is defined by:

$$d_{\mu,m} : x \in \mathbb{X} \mapsto \sqrt{\frac{1}{m} \int_0^m \delta_{\mu,l}(x)^2 dl}.$$

Proposition 3.12 Given a measure μ on a metric space \mathbb{X} and $x \in \mathbb{X}$, the application $m \in [0, 1[\mapsto d_{\mu,m}(x)$ is continuous on $[0, 1[$.

Proof: On $]0, 1[$, $m \mapsto \int_0^m \delta_{\mu,l}(x)^2 dl$ is continuous as the integral of a well-defined function. Furthermore, $m \mapsto \frac{1}{m}$ is continuous and thus $d_{\mu,m}$ is continuous. The only problem is for $m = 0$ due to the term $\frac{1}{m}$. It suffices to show the continuity in 0 and it will also guarantee that $d_{\mu,m}(x)$ is well-defined.

Let $m > 0$. Then by definition:

$$\forall 0 \leq l \leq m, \delta_{\mu,0}(x) \leq \delta_{\mu,l}(x) \leq \delta_{\mu,m}(x)$$

which translates for $d_{\mu,m}$:

$$\sqrt{\frac{1}{m} \int_0^m \delta_{\mu,0}(x)^2 dl} \leq \sqrt{\frac{1}{m} \int_0^m \delta_{\mu,l}(x)^2 dl} \leq \sqrt{\frac{1}{m} \int_0^m \delta_{\mu,m}(x)^2 dl}$$

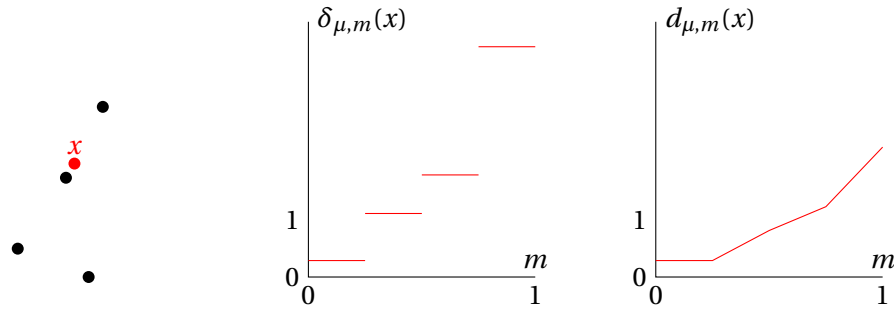
$$\delta_{\mu,0}(x) \leq d_{\mu,m}(x) \leq \delta_{\mu,m}(x)$$

By definition of the pseudo-distance, for any $\epsilon > 0$, $\mu(\bar{B}(x, \delta_{\mu,0} + \epsilon)) > 0$. Thus for any $0 < m < \mu(\bar{B}(x, \delta_{\mu,0} + \epsilon))$, $\delta_{\mu,m}(x) \leq \delta_{\mu,0} + \epsilon$. Consequently, the function $m \mapsto \delta_{\mu,m}(x)$ is continuous at 0. Hence $m \mapsto d_{\mu,m}(x)$ is continuous at 0 and $d_{\mu,0}(x) = \delta_{\mu,0}(x)$. ■

An interesting case for practical application is when the measure has a finite support and more precisely when μ is an empirical measure on a point set P . Figure 3.3 gives one example of a discontinuous pseudo-distance and continuous distance to a measure.

In this case, consider the parameter $k = mn$ where n is the number of points in P . Given a point $x \in \mathbb{X}$, the pseudo-distance $\delta_{\mu,m}(x)$ corresponds to the distance to the $[k]^{th}$ nearest neighbour of x in P . The distance to μ , $d_{\mu,m}(x)$ is the square average of the distance to the $[k]$ nearest neighbours of x . Let us write the points of P , $(p_1(x), \dots, p_n(x))$ such that for all $i < j$, $d_{\mathbb{X}}(p_i(x), x) \leq d_{\mathbb{X}}(p_j(x), x)$.

Proposition 3.13 Let P be a point cloud in a metric space \mathbb{X} , $m \in [0, 1[$ a mass parameter and


 Figure 3.3 – Empirical measure on a point cloud and values of $\delta_{\mu,m}$ and $d_{\mu,m}$

$x \in \mathbb{X}$. If μ is the empirical measure on P and $k = mn$ then

$$d_{\mu,m}(x) = \sqrt{\frac{1}{k} \left(\sum_{i=1}^{\lfloor k \rfloor} d_{\mathbb{X}}(p_i(x), x)^2 + (k - \lfloor k \rfloor) d_{\mathbb{X}}(p_{\lceil k \rceil}(x), x)^2 \right)}$$

Proof: Let us consider a ball of centre x and radius r . Fixing $d_{\mathbb{X}}(p_0(x), x) = 0$ and $d_{\mathbb{X}}(p_{n+1}(x), x) = \infty$, there exists an integer i such that $d_{\mathbb{X}}(p_i(x), x) \leq r < d_{\mathbb{X}}(p_{i+1}(x), x)$. Then $\mu(\bar{B}(x, r)) = \frac{i}{n}$ as the ball contains exactly i points from P . This implies that for any $i \leq n$ and any $\frac{i-1}{n} < m \leq \frac{i}{n}$, $\delta_{\mu,m}(x) = d_{\mathbb{X}}(p_i(x), x)$.

$$\begin{aligned} d_{\mu,m}(x)^2 &= \frac{1}{m} \int_0^m \delta_{\mu,l}(x)^2 dl \\ &= \frac{1}{m} \left(\sum_{i=1}^{\lfloor mn \rfloor} \int_{\frac{i-1}{n}}^{\frac{i}{n}} \delta_{\mu,l}(x)^2 dl + \int_{\frac{\lfloor mn \rfloor}{n}}^{\frac{mn}{n}} \delta_{\mu,l}(x)^2 dl \right) \\ &= \frac{1}{m} \left(\sum_{i=1}^k \frac{1}{n} d_{\mathbb{X}}(p_i(x), x)^2 + \frac{(k - \lfloor k \rfloor)}{n} d_{\mathbb{X}}(p_{\lceil k \rceil}(x), x)^2 \right) \\ &= \frac{1}{k} \left(\sum_{i=1}^k d_{\mathbb{X}}(p_i(x), x)^2 + d_{\mathbb{X}}(p_{\lceil k \rceil}(x), x)^2 \right) \end{aligned}$$

■

Remark that if k is an integer then the expression in Proposition 3.13 can be simplified as the term $k - \lfloor k \rfloor$ is zero. In practice, one always chooses k to be an integer which results in a simpler expression. From now on, unless specified otherwise, k will always be assumed to be an integer in order to simplify the proofs. However, all results adapt if k is not an integer.

3.4 Stability

In this section, we will study the stability of the distance to a measure with respect to the Wasserstein distance. It is a simple adaptation of the analysis in the Euclidean case [25].

First, let us introduce the notion of submeasure. A measure ν is a *submeasure* of a measure μ if for every $B \in \mathfrak{B}(\mathbb{X})$, $\nu(B) \leq \mu(B)$. Let $\text{Sub}_m(\mu)$ be the set of all submeasures of μ , that have a

total mass m . Then the distance to a measure μ at $x \in \mathbb{X}$ can be expressed as the Wasserstein distance between two measures, the Dirac mass $m\delta_x$ at x and a submeasure of μ of mass m . Intuitively, $d_{\mu,m}(x)$ is the minimal cost of moving a mass m from μ to x .

Proposition 3.14 *Let μ be a probability measure on a metric space \mathbb{X} , and let $m \in]0, 1[$ be a mass parameter. Then,*

$$d_{\mu,m}(x) = \min_{\nu \in \text{Sub}_m(\mu)} \frac{1}{\sqrt{m}} W_2(m\delta_x, \nu).$$

Given $x \in \mathbb{X}$ and $m > 0$, let $\mathcal{R}_{\mu,m}(x)$ be the set of the submeasures of μ with total mass m whose support is contained in the closed ball $\bar{B}(x, \delta_{\mu,m}(x))$ and whose restriction to the open ball $B(x, \delta_{\mu,m}(x))$ coincides with μ . The proof shows that $\mathcal{R}_{\mu,m}(x)$ is exactly the set of minimizers of Proposition 3.14.

In order to prove this theorem we need to introduce a few definitions. The *cumulative function* $F_\nu : \mathbb{R}^+ \rightarrow \mathbb{R}$ of a measure ν on \mathbb{R}^+ is the non-decreasing function defined by $F_\nu(y) = \nu([0, y])$. Its *generalized inverse* $F_\nu^{-1} : m \mapsto \inf\{t \in \mathbb{R} \mid F_\nu(t) > m\}$ is left-continuous.

Proof: If ν is a measure of total mass m on \mathbb{X} then there exists only one transport plan between ν and the Dirac mass $m\delta_x$. It transports every point of \mathbb{X} to x . Hence we get

$$W_2(m\delta_x, \nu)^2 = \int_{\mathbb{X}} d_{\mathbb{X}}(h, x)^2 d\nu(h).$$

Let $d_x : \mathbb{X} \rightarrow \mathbb{R}$ denote the distance function to the point x and let ν_x be the push-forward of ν by d_x . That is, for any subset I of \mathbb{R} , $\nu_x(I) = \nu(d_x^{-1}(I))$. Note that $F_{\nu_x}^{-1}(m) = \delta_{\nu,m}(x)$. Using the change of variable formula and the definition of the cumulative function, we get:

$$\int_{\mathbb{X}} d_{\mathbb{X}}(h, x)^2 d\nu(h) = \int_{\mathbb{R}^+} t^2 d\nu_x(t) = \int_0^m F_{\nu_x}^{-1}(l)^2 dl.$$

Suppose further that ν is a submeasure of μ , then $F_{\nu_x}(t) \leq F_{\mu_x}(t)$ for all $t > 0$. So, $F_{\nu_x}^{-1}(l) \geq F_{\mu_x}^{-1}(l)$ for all $l > 0$, and thus,

$$W_2(m\delta_x, \nu)^2 \geq \int_0^m F_{\mu_x}^{-1}(l)^2 dl = \int_0^m \delta_{\mu,l}(x)^2 dl = m d_{\mu,m}(x)^2.$$

This inequality implies that $d_{\mu,m}(x)$ is smaller than $\frac{1}{\sqrt{m}} W_2(m\delta_x, \nu)$ for any $\nu \in \text{Sub}_m(\mu)$.

Consider the case when the inequality is tight. Such a case happens when for almost every $l \leq m$, $F_{\nu_x}^{-1}(l) = F_{\mu_x}^{-1}(l)$. Since these functions are increasing and left-continuous, equality must hold for every such l . By the definition of the pushforward, this implies that $\nu(\bar{B}(x, \delta_{\mu,m}(x))) = m$, i.e., all the mass of ν is contained in the closed ball $\bar{B}(x, \delta_{\mu,m}(x))$, and that $\nu(B(x, \delta_{x,\mu}(m))) = \mu(B(x, \delta_{x,\mu}(m)))$. Because ν is a submeasure of μ this is true if and only if ν is in the set $\mathcal{R}_{\mu,m}(x)$ described before the proof. Thus $\mathcal{R}_{\mu,m}(x)$ is exactly the set of submeasures $\nu \in \text{Sub}_m(\mu)$ such that $d_{\mu,m}(x) = \frac{1}{\sqrt{m}} W_2(m\delta_x, \nu)$.

To conclude the proof we only need to show that there exists at least one measure $\mu_{x,m}$ in the set $\mathcal{R}_{\mu,m}(x)$. If $\mu(\bar{B}(x, \delta_{\mu,m}(x))) = m$, then $\mu_{x,m} = \mu|_{\bar{B}(x, \delta_{\mu,m}(x))}$ is an obvious choice. The only difficulty is when the boundary $\partial B(x, \delta_{\mu,m}(x))$ of the ball has too much mass. In this case we

uniformly rescale the mass contained in the bounding sphere such that the measure $\mu_{x,m}$ has total mass m . More precisely we let:

$$\mu_{x,m} = \mu|_{B(x,\delta_{\mu,m}(x))} + (m - \mu(B(x,\delta_{\mu,m}(x)))) \frac{\mu|_{\partial B(x,\delta_{\mu,m}(x))}}{\mu(\partial B(x,\delta_{\mu,m}(x)))}.$$

We hence have $\frac{1}{\sqrt{m}} W_2(m\delta_x, \mu_{x,m}) = d_{\mu,m}(x)$. ■

From this result, we have the following Wasserstein stability guarantee for the distance to a measure.

Theorem 3.15 *Let μ and ν be two probability measures on a metric space \mathbb{X} and let $m \in]0, 1[$ be a mass parameter. Then:*

$$\|d_{\mu,m} - d_{\nu,m}\|_{\infty} \leq \frac{1}{\sqrt{m}} W_2(\mu, \nu).$$

Proof: Using Proposition 3.14, we get that $\sqrt{m} d_{\mu,m}(x) = W_2(m\delta_x, \mu_{x,m})$, where $\mu_{x,m} \in \mathcal{R}_{\mu,m}(x)$. Let π be an optimal transport plan between μ and ν , i.e., a transport plan between μ and ν such that

$$\int_{\mathbb{X} \times \mathbb{X}} d_{\mathbb{X}}(x, y)^2 d\pi(x, y) = W_2(\mu, \nu)^2.$$

Let us consider the submeasure $\mu_{x,m}$ of μ . There exists $\tilde{\pi}$ a submeasure of π that transports $\mu_{x,m}$ to a submeasure $\tilde{\nu}$ of ν of mass m and

$$W_2(\mu_{x,m}, \tilde{\nu}) \leq W_2(\mu, \nu).$$

Using Proposition 3.14 again, we get for any $x \in \mathbb{X}$, $\sqrt{m} d_{\nu,m}(x) \leq W_2(m\delta_x, \tilde{\nu})$. Thus,

$$\begin{aligned} \sqrt{m} d_{\nu,m}(x) &\leq W_2(m\delta_x, \tilde{\nu}) \leq W_2(m\delta_x, \mu_{x,m}) + W_2(\tilde{\nu}, \mu_{x,m}) \\ &\leq \sqrt{m} d_{\mu,m}(x) + W_2(\mu, \nu). \end{aligned}$$

The roles of μ and ν can be reversed to conclude the proof. ■

The combination of Proposition 3.14 and Theorem 3.15 implies that the function $x \mapsto d_{\mu,m}(x)$ is 1-Lipschitz.

Corollary 3.16 *Let μ be a probability measure on a metric space \mathbb{X} and let $m \in [0, 1[$ be a mass parameter. The function $x \mapsto d_{\mu,m}(x)$ is 1-Lipschitz.*

Proof: Remark that if $m = 0$ then $d_{\mu,0}$ is the distance to $\text{Supp}(\mu)$ which is 1-Lipschitz. Assume that $m > 0$ and let x and y be two points of \mathbb{X} . Proposition 3.14 implies

$$d_{\mu,m}(x) = \min_{\nu \in \text{Sub}_m(\mu)} \frac{1}{\sqrt{m}} W_2(m\delta_x, \nu).$$

Let $\nu \in \text{Sub}_m(\mu)$. We have:

$$\begin{aligned} d_{\mu,m}(x) &\leq \frac{1}{\sqrt{m}} W_2(m\delta_x, \nu) \\ &\leq \frac{1}{\sqrt{m}} W_2(m\delta_y, \nu) + \frac{1}{\sqrt{m}} W_2(m\delta_x, m\delta_y) \\ &= \frac{1}{\sqrt{m}} W_2(m\delta_y, \nu) + d_{\mathbb{X}}(x, y) \end{aligned}$$

Thus $d_{\mu,m}(x) - d_{\mathbb{X}}(x, y) \leq d_{\mu,m}(y)$. The roles of x and y can be interchanged to obtain us the 1-Lipschitz property as $|d_{\mu,m}(x) - d_{\mu,m}(y)| \leq d_{\mathbb{X}}(x, y)$. ■

3.5 Discriminating results

We just showed that two measures that are close with respect to the Wasserstein distance give close distance to measure functions. A natural question arises. Is the opposite true? If two measures gives distance to measure functions that are close to each other, what can we say about the measures themselves? In this section and the next, we provide some partial answers to this question.

We will first show that if two distances to measure are equal everywhere and for all values of the mass parameter on a Euclidean space then the two measures are equal. Then we relax some hypotheses on the mass or the number of points where the distances are known while restricting the set of measures to the ones with finite support. All presented results are qualitative and finding quantitative results is an ongoing research direction.

3.5.1 Correspondence between measure and distance to a measure

The distances to measure and measures have a one-to-one correspondence in the sense that a distance to a measure defined for all points of \mathbb{R}^d and all masses m can be generated by only one probability measure on \mathbb{R}^d . This translates in the following theorem:

Theorem 3.17 *Let μ and ν be two probability measures on \mathbb{R}^d , then:*

$$\left(\forall x \in \mathbb{R}^d, \forall m \in [0, 1[, d_{\mu,m}(x) = d_{\nu,m}(x) \right) \Leftrightarrow \mu = \nu$$

Before proving this result, let us recall an elementary lemma from measure theory.

Lemma 3.18 ([87, Theorem 1.19]) *Let A_n be a family of Borel sets such that $A_{i+1} \subset A_i$ for all i and let μ be a non-negative measure such that $\mu(A_1)$ is finite. Writing $A = \bigcap_{n=1}^{\infty} A_n$:*

$$\lim_{n \rightarrow \infty} \mu(A_n) = \mu(A)$$

Let us now prove Theorem 3.17.

Proof: The construction of the distance to a measure is unique, hence if $\mu = \nu$, then $d_{\mu,m} = d_{\nu,m}$. We need to prove that if $d_{\mu,m}(x) = d_{\nu,m}(x)$ for all x and m , then the two measures μ and ν are equal.

Assume that it is not the case. Then there exists a Borel subset $U \in \mathfrak{B}(\mathbb{R}^d)$ such that $\mu(U) \neq \nu(U)$. $\mathfrak{B}(\mathbb{R}^d)$ is generated by the balls of \mathbb{R}^d . Thus, there exists a point x and a radius r such that:

$$\mu(B(x, r)) \neq \nu(B(x, r))$$

Without loss of generality, let us assume that $m_0 = \mu(B(x, r)) < \nu(B(x, r)) = m_2$. Let $A_n = B(x, r) \setminus \bar{B}(x, r - \frac{1}{n})$ for $n \geq n_0 = \lceil \frac{1}{r} \rceil$. Then for all n , $A_{n+1} \subset A_n$ and $A = \cap_{n=0}^{\infty} A_n$. Applying Lemma 3.18:

$$\exists n, \nu(A_n) < m_2 - m_0$$

This implies:

$$\exists \epsilon > 0, m_0 < \nu(B(x, r - \epsilon)) = m_1 \leq m_2$$

By definition of the pseudo-distance to a measure:

$$\forall m, m_0 < m < m_1, \delta_{\mu, m}(x) \geq r > r - \epsilon \geq \delta_{\nu, m}(x)$$

Our initial assumption states in particular that $d_{\mu, m_0}(x) = d_{\nu, m_0}(x)$.

$$\begin{aligned} d_{\nu, m_1}(x)^2 &= \frac{1}{m_1} \int_0^{m_1} \delta_{\nu, m}^2 dm \\ &= \frac{1}{m_1} \int_0^{m_0} \delta_{\nu, m}^2 dm + \frac{1}{m_1} \int_{m_0}^{m_1} \delta_{\nu, m}^2 dm \\ &\leq \frac{m_0}{m_1} d_{\nu, m}(x)^2 + \frac{1}{m_1} (r - \epsilon)^2 dm \\ &< \frac{m_0}{m_1} d_{\mu, m}(x)^2 + \frac{1}{m_1} \int_{m_0}^{m_1} \delta_{\mu, x}(x)^2 dm \\ &= d_{\mu, m_1}(x)^2 \end{aligned}$$

Thus we found an x and an m such that $d_{\nu, m}(x) \neq d_{\mu, m}(x)$, which contradict the assumption. Hence $\mu = \nu$. ■

3.5.2 Recovering a measure from its distance

The correspondence between distance to a measure functions and measures makes it conceivable to reconstruct a measure knowing the distance to a measure function. This means that for a given Borel subset, we can compute its measure. In this section, we introduce some partial inversion for the pseudo distance $\delta_{\mu, m}$ to a measure μ . As we explained in section 3.3, the pseudo-distance is not continuous. We then need a right and a left inverse. For commodity we write them:

$$\begin{aligned} m_{\mu, x}^-(r) &= \sup\{m | \delta_{\mu, m}(x) < r\} \\ m_{\mu, x}^+(r) &= \inf\{m | \delta_{\mu, m}(x) > r\} \end{aligned}$$

These inverses allow us to compute the mass of balls centred in x and with radius r thanks to the following technical lemma:

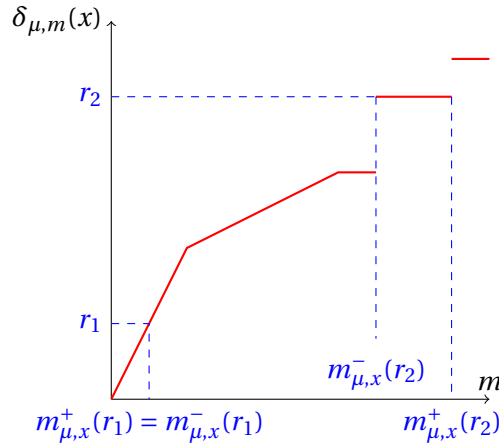


Figure 3.4 – Left and right inverses of the pseudo-distance $\delta_{\mu,m}(x)$

Lemma 3.19 Let μ be a probability measure on a metric space \mathbb{X} and $x \in \mathbb{X}$. For any $r \geq 0$,

$$m_{\mu,x}^-(r) = \mu(B(x, r))$$

$$m_{\mu,x}^+(r) = \mu(\bar{B}(x, r))$$

Proof: Let us fix $r_0 \geq 0$ and show that $m_{\mu,x}^-(r_0) = \mu(B(x, r_0))$.

Let m be such that $\delta_{\mu,m}(x) < r_0$. There exists r_1 such that $\delta_{\mu,m}(x) < r_1 < r_0$. Thus $\inf\{r \mid \mu(\bar{B}(x, r)) > m\} < r_1 < r_0$. Hence $\mu(B(x, r_0)) \geq \mu(\bar{B}(x, r_1)) > m$. By definition, $m_{\mu,x}^-(r_0) \leq \mu(B(x, r_0))$.

Now assume that $\mu(B(x, r_0)) > m_{\mu,x}^-(r_0)$. Remark that $\mu(B(x, r_0))$ is the limit of $\mu(\bar{B}(x, r))$ when r increases and tends to r_0 . Hence, there exists $r_1 < r_0$ such that $\mu(\bar{B}(x, r_1)) > m_{\mu,x}^-(r_0)$. Then there exists m_2 such that $m_{\mu,x}^-(r_0) < m_2 < \mu(\bar{B}(x, r_1))$, which means that $\delta_{\mu,m_2}(x) \leq r_1$. Finally $m_{\mu,x}^-(r_0) \geq m_2 > m_{\mu,x}^-(r_0)$ which gives a contradiction. Thus $m_{\mu,x}^-(r_0) = \mu(B(x, r_0))$.

Let us now prove $m_{\mu,x}^+(r_0) = \mu(\bar{B}(x, r_0))$.

Fix $m > m_{\mu,x}^+(r_0)$. Then $\delta_{\mu,m}(x) > r_0$, id est $\inf\{r \mid \mu(\bar{B}(x, r)) > m\} > r_0$. Hence, $\mu(\bar{B}(x, r_0)) \leq m$ and by definition $\mu(\bar{B}(x, r_0)) \leq m_{\mu,x}^+(r_0)$.

Now assume that $\mu(\bar{B}(x, r_0)) < m_{\mu,x}^+(r_0)$. Remark that $\mu(\bar{B}(x, r_0))$ is the limit of $\mu(\bar{B}(x, r))$ when r decreases and tends to r_0 . Hence, there exists $r > r_0$ and m such that $\mu(\bar{B}(x, r)) < m < m_{\mu,x}^+(r_0)$. Thus $\delta_{\mu,m}(x) \geq r > r_0$ which implies $m_{\mu,x}^+(r_0) \leq m < m_{\mu,x}^+(r_0)$. This is a contradiction and then $m_{\mu,x}^+(r_0) = \mu(\bar{B}(x, r_0))$. ■

Given a measure μ and a point x , the function $m \mapsto d_{\mu,m}(x)^2$ defined on $[0, 1[$ is non-decreasing, continuous and differentiable. In fact, it is given by the relation

$$d_{\mu,m}^2(x) = \frac{1}{m} \int_0^m \delta_{\mu,l}(x)^2 dl$$

where $l \mapsto \delta_{\mu,l}(x)$ is a non-decreasing function. Knowing $d_{\mu,m}(x)$ for any mass m , it is possible to compute the value of $\delta_{\mu,m}(x)$ for any value of $m \in [0, 1[$. Moreover, we can compute $m_{\mu,x}^-(r)$ and $m_{\mu,x}^+(r)$.

3.5.3 Application to the real line

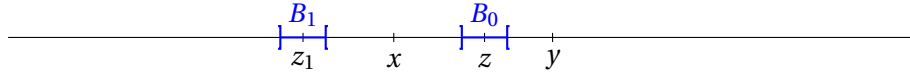
Using this in dimension 1, we design an algorithm to reconstruct a measure μ knowing the value of $d_{\mu,m}$ at two points for every $m \in [0, 1[$. In fact, we show that we recover the measure of a ball with an arbitrary precision ϵ , by which we mean that difference between the computed value and the real one is at most ϵ .

Theorem 3.20 *Let μ be a probability measure on \mathbb{R} . Given two distinct points x and y in \mathbb{R} and $d_{\mu,m}(x)$ and $d_{\mu,m}(y)$ for all $m \in [0, 1[$, there exists an algorithm that computes $\mu(B(z, \rho))$ with precision ϵ for any $z \in \mathbb{R}$, $\rho > 0$ and $\epsilon > 0$. Moreover, if μ has compact support, then $\mu(B(z, \rho))$ can be exactly computed in finite time.*

Proof: Remark first that, if $z = x$, then $\mu(B(z, \rho)) = m_{\mu,x}^-(\rho)$. Hence, let us assume that $z \neq x$ and $z \neq y$. The principle of the algorithm is to compute the mass of the ball $B_0 = \mu(z, \rho)$ by iteratively adding and removing masses from a sequence of intervals $(B_i)_{i \geq 0}$ centred on points $(z_i)_{i \geq 0}$. At each step we compute a quantity s_i that will be close to $\mu(B_0)$ for some i determined during the execution of the algorithm.

Construction:

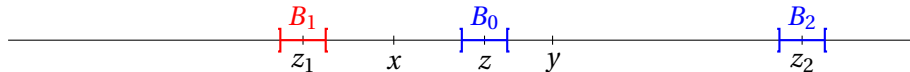
Without loss of generality, we assume that $x < y$. By convention we write $z_0 = z$. First we consider the case $z > x$ and $z - x \geq \rho$. We build z_1 to be the symmetric of z_0 by x and fix $B_1 = B(z_1, \rho) =]z_1 - \rho, z_1 + \rho[$. Note that $B_0 \cap B_1 = \emptyset$.



Remark that $a_1 = \mu(B_0) + \mu(B_1)$ can be computed because $a_1 = \mu(B(x, z - x + \rho)) - \mu(\bar{B}(x, z - x - \rho))$ and using Lemma 3.19

$$a_1 = m_{\mu,x}^-(z + \rho - x) - m_{\mu,x}^+(z - \rho - x) = s_1.$$

We use information at point y . We define z_2 to be the symmetric of z_1 by y and fix $B_2 = B(z_2, \rho)$.



Again the mass $\mu(B_1) + \mu(B_2)$ can be computed using Lemma 3.19.

$$a_2 = m_{\mu,y}^-(z_2 - y + \rho) - m_{\mu,y}^+(z_2 - y - \rho)$$

Remark that B_2 is disjoint from B_1 and B_0 . Moreover, $s_2 = a_1 - a_2 = \mu(B_0) - \mu(B_2)$. We carry on the iterative construction of z_{i+1} by taking the symmetric of z_i by x if i is even and by y if i is odd. B_i is defined as the open ball of centre z_i and radius ρ , which is disjoint from all B_j when $j < i$, while $s_{i+1} = s_i + (-1)^i a_{i+1}$ where $a_{i+1} = \mu(B_i) + \mu(B_{i+1})$. Every a_i is computable as the difference of mass between two balls centred either in x or in y .

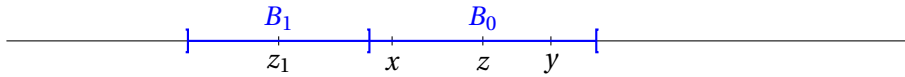
$$a_{2i+1} = \mu(B(x, z_{2i+1} + \rho)) - \mu(\bar{B}(x, z_{2i+1} - \rho)) = m_{\mu,x}^-(z_{2i+1} - x + \rho) - m_{\mu,x}^+(z_{2i+1} - x - \rho)$$

$$a_{2i} = \mu(B(y, z_{2i} + \rho)) - \mu(\bar{B}(y, z_{2i} - \rho)) = m_{\mu,y}^-(y + z_{2i} + \rho) - m_{\mu,y}^+(y + z_{2i} - \rho)$$

By a trivial recursion, for any $n > 0$,

$$s_n = \sum_{i=1}^n (-1)^{i+1} a_i = \mu(B_0) + (-1)^{n+1} \mu(B_n).$$

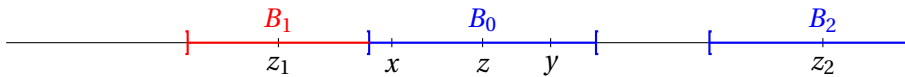
If $z < y$ and $z - y \geq \rho$, then the symmetric construction obtained by reversing the roles of x and y works to obtain the same final relation. However, if $z > x$, $z - x < \rho$ and $y - x < 2\rho$, this construction is not possible because B_0 and B_1 will no longer be disjoint. We need to slightly change the construction of the sequence B_i . Be careful that the sets will no longer be open intervals and it will change the way the various a_i are computed. First, we build B_1 such that $B_1 \cup B_0$ is an open ball centred in x and $B_1 \cap B_0 = \emptyset$. This means that $B_1 =]2x - z - \rho, z - \rho]$ is closed on the right side. To come back to the first construction, we fix $z_1 = x - \rho$ to be the mean of B_1 and $\rho' = z - x$ to be half its diameter. Hence $B_1 =]z_1 - \rho', z_1 + \rho']$.



Keeping the definition $a_i = \mu(B_i) + \mu(B_{i+1})$, we have $a_1 = \mu(B(x, z + \rho - x)) = m_{\mu,x}^-(z + \rho - x)$. We build iteratively B_{i+1} as before, taking the symmetric of B_i either by x or y , depending on the parity of i . The definitions of z_{i+1} , a_{i+1} and s_{i+1} follow naturally. Notice that B_0 and B_2 are not necessarily disjoint. However, the construction is the same as before starting with B_1 and for any $i > j > 0$, $B_i \cap B_j = \emptyset$. Be also careful that the computation of a_{i+1} is slightly different from before, the function m^+ being replaced by m^- .

$$a_{2i+1} = \mu(B(x, z_{2i+1} + \rho)) - \mu(B(x, z_{2i+1} - \rho)) = m_{\mu,x}^-(z_{2i+1} - x + \rho) - m_{\mu,x}^-(z_{2i+1} - x - \rho)$$

$$a_{2i} = \mu(B(y, z_{2i} + \rho)) - \mu(B(y, z_{2i} - \rho)) = m_{\mu,y}^-(y + z_{2i} + \rho) - m_{\mu,y}^-(y + z_{2i} - \rho)$$



Termination of the algorithm:

Remark that $\mu(B_0)$ is bounded by the values (s_i) as for any $i \geq 1$:

$$s_{2i} \leq \mu(B_0) \leq s_{2i+1}.$$

Given the value ϵ , we can decide to stop the algorithm as soon as $|s_{i+1} - s_i| < \epsilon$. Then $|s_i - \mu(B_0)| < \epsilon$. We need to show that there exists such a i .

Fix n such that $n\epsilon > 1$. The measure μ is a probability measure and thus is positive with total mass 1. Moreover, all the B_i for $i \geq 1$ are pairwise disjoint and hence there exists $j \in [1, n]$ such that

$$a_{2j} = \mu(B_{2j-1}) + \mu(B_{2j}) < \epsilon.$$

Otherwise:

$$\mu(\mathbb{R}) \geq \sum_{i=1}^{2n} \mu(B_i) = \sum_{j=1}^n a_{2j+1} > n\epsilon > 1$$

When the measure has compact support then there exists a N such that for any $n \geq N$, $B_n \cap \text{Supp}(\mu) = \emptyset$ because (z_{2i}) and (z_{2i+1}) are two unbounded monotonic sequence and hence $\lim_{n \rightarrow \infty} |z_n| = +\infty$. Thus for any $n \geq N$, we have $s_n = s_{n+1} = \mu(B_0)$. ■

3.5.4 Higher dimensional reconstruction for finite support measures

Unfortunately, the algorithm in dimension 1 can not be transposed in higher dimensions. However, by considering only measures with finite support, we can reconstruct exactly a measure μ in \mathbb{R}^d by knowing the value of $d_{\mu,m}$ at $d+2$ points for all masses $m \in [0, 1[$.

HD-RECONSTRUCTION

1. Query $d_{\mu,m}(x_i)$ for $d+1$ points (x_1, \dots, x_{d+1}) in general position.
 2. For every i , compute the set $W_i = \{r \mid m_{\mu,x_i}^-(r) \neq m_{\mu,x_i}^+(r)\}$.
 3. Compute a superset Ω of $\text{Supp}(\mu)$ by taking $\Omega = \cup_{\sigma \in \prod_{i=0}^d W_i} \cap_{j=1}^d S(x_j, \sigma_j)$, where σ_j is the j^{th} component of σ and $S(x, r)$ is the sphere of centre x and radius r .
 4. Assign correct masses to the points of Ω using the relation $\mu(S(x_i, r)) = m_{\mu,x_i}^+(r) - m_{\mu,x_i}^-(r)$, making a query in another point if necessary.
-

Theorem 3.21 *Let μ be a probability measure on \mathbb{R}^d with finite support. Assuming that we can query the function $m \mapsto d_{\mu,m}(x)$ for any x , the algorithm HD – reconstruction reconstructs the measure μ with $d+2$ queries.*

Proof: The remarks made in section 3.5.2 makes it possible to compute the set W_i for all i . Given a point x_i , it is sufficient to derive the function $m \mapsto d_{\mu,m}(x_i)$ and look for its discontinuities. This way we can build the set $W = \prod_{i=0}^d W_i$.

The computation of the set Ω can be reduced to solving a set of linear systems. We lift our space on the paraboloid P embedded in \mathbb{R}^{d+1} such that $y_{d+1}^2 = \sum_{i=1}^d y_i^2$ for any point $y = (y_1, \dots, y_d) \in \mathbb{R}^d$. A sphere of \mathbb{R}^d becomes a hyperplane in \mathbb{R}^{d+1} . Given $\sigma \in \prod_{i=0}^d W_i$, the set $\cap_{j=1}^d S(x_j, \sigma_j)$ is the solution of a linear system of $d+1$ equations expressing the intersection of $d+1$ hyperplanes. Due to the assumption on the general position of (x_1, \dots, x_{d+1}) , this system is non-degenerate and has a unique solution. We define Ω_0 to be the set of all points that are solution of at least one such linear system. We are only interested in the points located on the paraboloid P . Hence, we take the intersection of Ω_0 and P before projecting back on \mathbb{R}^d to obtain the set Ω .

Given a point x and a radius r , Lemma 3.19 gives $\mu(S(x, r)) = \mu(\bar{B}(x, r)) - \mu(B(x, r)) = m_{\mu,x}^+(r) - m_{\mu,x}^-(r)$. Thus, we have $\text{Supp}(\mu) \subset \Omega$. Moreover, it gives a system of $\sum_{i=1}^d |W_i|$ equations on

the mass of points. For any $i \in [1, d]$ and $r \in W_i$, the mass of $\mu(S(x_i, r) \cap \Omega)$ is known. If this intersection is a single point then we assign the corresponding mass to the point. This is always possible because μ is a valid assignment for the system. However, it is possible that multiple assignments exist.

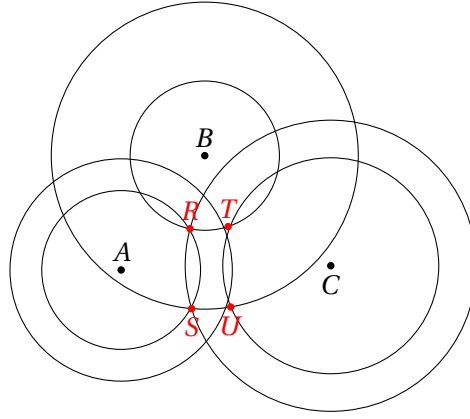


Figure 3.5 – Case where masses can not be identified

This happens where it is impossible to distinguish between points as shown in Figure 3.5. Ω is equal to the red points (R, S, T, U) and we know α, β, γ such that:

$$\mu(R) + \mu(S) = \alpha = 1 - (\mu(T) + \mu(U))$$

$$\mu(R) + \mu(T) = \beta = 1 - (\mu(S) + \mu(U))$$

$$\mu(T) + \mu(U) = \gamma = 1 - (\mu(R) + \mu(S))$$

This system is not well conditioned. For example, if $\alpha = \beta = \gamma = \frac{1}{2}$ then $\mu(R) = \mu(S) = \mu(T) = \mu(U) = \frac{1}{4}$ is a solution but $\mu(R) = \mu(U) = \frac{1}{2}$ is another one. Notice that the three points A, B and C are in general position. In this case, we can build a family of measures which are candidates for μ . To recover μ exactly, we need to add a new query point. The way to remove all ambiguities is to choose x_{d+2} such that no two points of Ω are at equal distance of x_{d+2} . In other words, x_{d+2} is outside the union of all bisector hyperplanes of points of Ω , which happens with probability 1 if the point is randomly selected.

Remark that the cardinality of $\text{Supp}(\mu)$ is trivially lower bounded by the height of W and that the necessity to use an additional point comes from very specific positions. ■

3.6 Relation to higher order Voronoi and power diagrams

In this section, we explore the relations between the distance to a measure and two families of diagrams. The relation to the k^{th} -order Voronoi diagram gives an interesting result on reconstructing a measure from its distance with results similar to those of Section 3.5. This is in turn related to power diagrams that provide the foundation upon which the approximation techniques of Chapter 4 are built.

3.6.1 k^{th} -order Voronoi diagrams

The cell of a Voronoi diagram is defined as the set of points that share the same nearest neighbour in P .

Definition 3.22 Let P be a point set in a metric space \mathbb{X} . For any $x \in P$, the Voronoi cell of x is the set:

$$V(x) = \{y \in \mathbb{X} \mid \forall p \in P, d_{\mathbb{X}}(x, y) \leq d_{\mathbb{X}}(y, p)\}$$

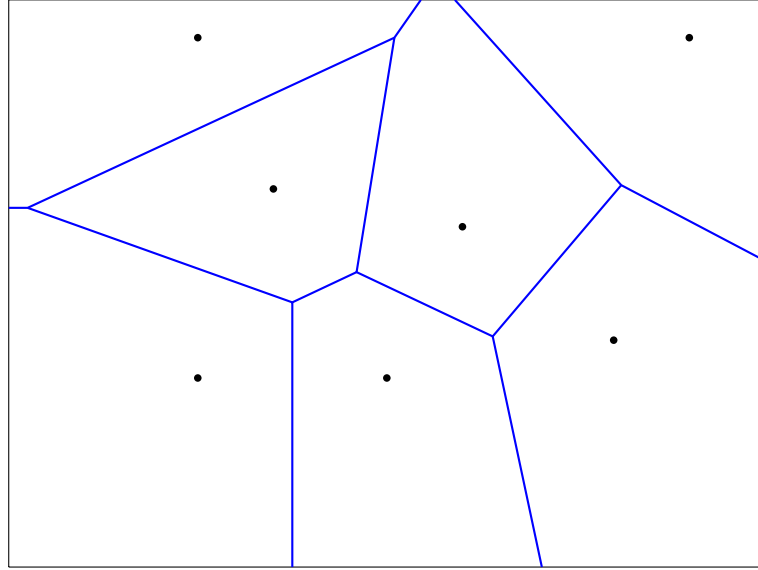


Figure 3.6 – A point cloud and its Voronoi diagram

The Voronoi diagram is the collection of Voronoi cells and their intersections. In the case of Euclidean spaces, it corresponds to the dual of the famous Delaunay triangulation. More details can be found in any classic book on computational geometry, for example [15]. This notion can be easily extended to higher order diagrams. The k^{th} -order Voronoi diagram is built by looking at the k nearest neighbours.

Definition 3.23 Let P be a point set in a metric space \mathbb{X} . For any subset (x_1, \dots, x_k) of k elements of P , the corresponding k^{th} -order Voronoi cell is the set:

$$V^k(x_1, \dots, x_k) = \{y \in \mathbb{X} \mid \forall p \in P \setminus (x_1, \dots, x_k), \forall i \in [1, k], d_{\mathbb{X}}(x_i, y) \leq d_{\mathbb{X}}(p, y)\}$$

Considering an empirical measure μ defined from a point set P , the relation between the distance to μ for a mass $m = \frac{k}{n}$ and the k^{th} -order Voronoi diagram is direct. All points in the same k^{th} -order Voronoi cell share the same k nearest neighbours and thus $d_{\mu, m}^2$ restricted to the cell is a quadratic form.

$$\forall x \in V^k(x_1, \dots, x_k), d_{\mu, m}(x)^2 = \frac{1}{k} \sum_{i=1}^k d_{\mathbb{X}}(x_i, x)^2$$

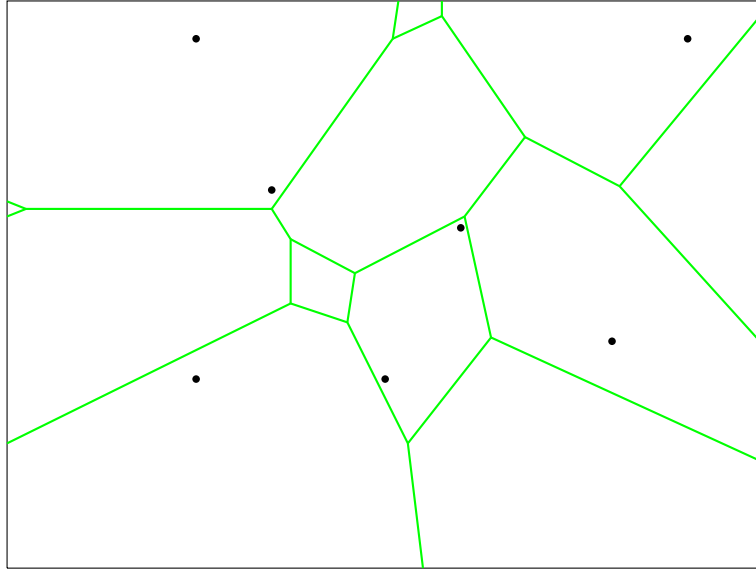


Figure 3.7 – The second order Voronoi diagram of the previous point cloud

This remark is especially interesting when the metric space \mathbb{X} is Euclidean. It allow us to identify a measure with finite support knowing only its distance for a sufficiently small mass m .

Theorem 3.24 *Let f be a function on a Euclidean space \mathbb{R}^d , $d \geq 2$, and $m \in [0, \frac{1}{2}]$. If there exists an empirical measure μ with finite support P in general position such that $f = d_{\mu,m}$ and for any $x \in \mathbb{R}^d$, $\mu(x) \leq m$, then the measure μ is uniquely defined.*

Proof: For the sake of simplicity, we assume that $k = mn$ is an integer. Due to the hypotheses, f is a quadratic form on each of the k^{th} -order Voronoi cells of P . Thus it is differentiable in the interior of all these cells. Let us define the set $\mathcal{D} = \{x \in \mathbb{R}^d \mid f \text{ is not differentiable at } x\}$. We first show that \mathcal{D} corresponds exactly to the boundary of k^{th} -order Voronoi cells and thus to the k^{th} -order Voronoi diagram.

Let x be a point on the boundary of a k^{th} -order Voronoi cell. Then there exists two sets (x_1, \dots, x_{k-1}, y) and (x_1, \dots, x_{k-1}, z) such that $x \in V^k(x_1, \dots, x_{k-1}, y) \cap V^k(x_1, \dots, x_{k-1}, z)$ and $y \neq z$. Moreover, we can choose y and z such that the two cells are not reduced to the point x . Let us assume that f is differentiable at x . Then f^2 is also differentiable at x and we can compute the gradient of f^2 .

$$\begin{aligned} \nabla f^2(x) &= \frac{2}{k} \left(\sum_{i=1}^{k-1} (x - x_i) + x - y \right) \\ &= \frac{2}{k} \left(\sum_{i=1}^{k-1} (x - x_i) + x - z \right) \end{aligned}$$

However, $y \neq z$ and thus we have a contradiction. The set \mathcal{D} corresponds exactly to the intersections of two or more k^{th} -order Voronoi cells.

The k^{th} -order Voronoi diagram is an affine diagram. Let us consider an edge e of this diagram. Then there exist d points $(x_1, \dots, x_d) \in P^d$ such that e is contained in the bisector hyperplanes of x_i and x_j for all $i \neq j$. We denote a and b the two extremities of e . a is the centre of the circumsphere to the d -dimensional simplex σ whose vertices are (x_1, \dots, x_d) and another point in P , denoted x_{d+1} . The $d - 1$ dimensional faces of the k^{th} -order Voronoi diagram containing a are included inside the bisector hyperplanes of pairs of edges of σ . Knowing these orthogonality conditions, the positions of the points (x_1, \dots, x_d) possesses one degree of freedom and σ is defined up to a homothety. Figure 3.8 illustrates the situation for $d=2$. The dashed lines correspond to the directions on which the points can be placed and the blue triangles represent some possible simplices σ .

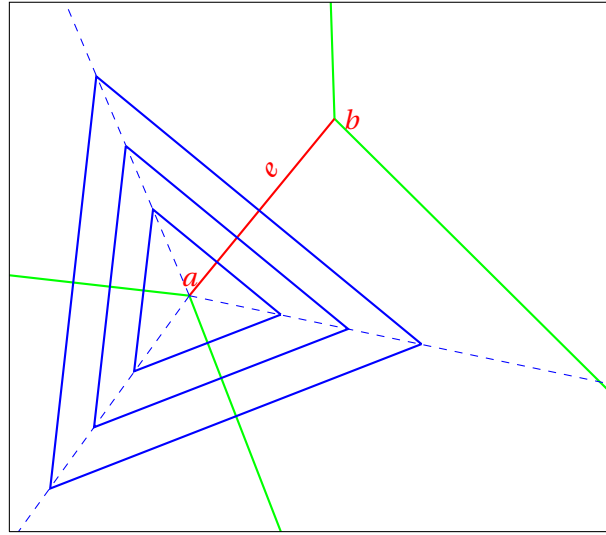


Figure 3.8 – An edge with similar triangles having A as orthocentre

Assuming that the edge e is non-degenerate, consider one of the full-dimensional cell incident to e . Only one point y from the set $\{x_1, \dots, x_d\}$ is part of the subset associated with the cell. The homotheties obtained for the simplices using the extremities a and b of e defines exactly y as the intersection of two lines. Consequently, all the vertices of σ are uniquely defined. Thus, any point p of P such that there exists a non-degenerate edge e with p in the set (x_1, \dots, x_d) can be reconstructed, which is equivalent to say that at least one k^{th} -order cell is not associated with p .

Having a mass parameter $m < \frac{1}{2}$ guarantees that, for all $p \in P$, there exists a cell that is not associated with p . Let u be a direction. The hyperplane H orthogonal to u and containing p subdivides the spaces in two. We choose u such that $H \cap P = \{p\}$. One of the half-space defined by h contains more than k points because $m < \frac{1}{2}$. Hence the unbounded k^{th} -order cell in the direction of u in this half-space is not associated with p . ■

The general position assumption is not needed as long as we can find a non-degenerate edge for every point p . The assumption is made here to avoid a more complicated discussion.

The hypothesis on the mass stating that $m < \frac{1}{2}$ is necessary as illustrated by Figure 3.9. In this example, the point located at the centre of the circle is part of all Voronoi cells for $m \geq \frac{1}{2}$ and thus cannot be precisely placed. In fact, if $m > \frac{1}{2}$, it will be possible to slightly move it without modifying the diagram. The centre of the circle is associated with all cells of the 10^{th} -order Voronoi diagram. One cell is highlighted as well as its associated points.

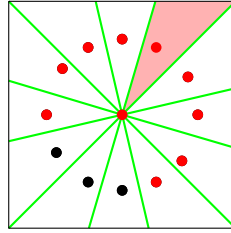


Figure 3.9 – A point cloud where the central point is associated with all cells for $k = 10$

If we work on the real line, the absence of edges makes it necessary to know the parameter m a priori in order to reconstruct the measure. However, assuming that m is known, we can compute n as the number of cells is exactly $n - k$ and the reconstruction is easy. Ordering the points of P such that $p_1 \leq p_i \leq p_n$, and the points of \mathcal{D} as $y_1 \leq y_j \leq y_{n-k-1}$, we have the relations:

$$y_j - p_j = p_{j+k} - y_j$$

This holds for all j and due to the assumption that $m < \frac{1}{2}$, the system is well conditioned and admits a unique solution.

The set of points P used to defined the empirical measure μ was implicitly assumed to not contain any point with multiplicity 2. This uniqueness result can be extended to measures that do not correspond to empirical measure or to the case of P having points with multiplicity. However, it is necessary to assume that no point of $\text{Supp}(\mu)$ accumulates more than a mass m and that there exists a cell not associated with p for all $p \in P$. Then the construction of $\text{Supp}(\mu)$ is identical as the one in the previous proof. The mass of each point p has to be computed at the end. Remark that the support points of the power distance $d_{\mu,m}$ can be found and that they are barycentres of points of $\text{Supp}(\mu)$ taken with their weights. The number of cells being more than the number of points, we can build a well-conditioned system of equations to reconstruct μ .

3.6.2 Power diagrams

k^{th} -order Voronoi diagrams are affine diagrams and thus are related to *power diagrams*. A power diagram is defined using a function called a *power distance*.

Definition 3.25 Let P be a compact subset of a metric space \mathbb{X} . Given a function $w : P \rightarrow \mathbb{R}$, the

power distance f associated with (P, w) is defined by:

$$f(x) = \inf_{p \in P} \left(\sqrt{d_{\mathbb{X}}(p, x)^2 + w(p)^2} \right)$$

P is called the support of f .

For the sake of simplicity, the weight $w(p)$ will be written as w_p in the rest of the dissertation.

Lemma 3.26 *A power distance f associated with (P, w) in a metric space \mathbb{X} is 1-Lipschitz.*

Proof: Consider x and y in \mathbb{X} . For any $\epsilon > 0$, there exists $p \in P$ such that $w_p^2 + d_{\mathbb{X}}(x, p)^2 \leq f(x)^2 + \epsilon^2$.

$$\begin{aligned} f(y)^2 &\leq w_p^2 + d_{\mathbb{X}}(y, p)^2 \\ &\leq w_p^2 + d_{\mathbb{X}}(x, p)^2 + d_{\mathbb{X}}(x, y)^2 + 2d_{\mathbb{X}}(x, p)d_{\mathbb{X}}(y, p) \\ &\leq f(x)^2 + d_{\mathbb{X}}(x, y)^2 + 2f(x)d_{\mathbb{X}}(y, p) + \epsilon(\epsilon + 2d_{\mathbb{X}}(y, p)) \\ &\leq (f(x) + d_{\mathbb{X}}(x, y))^2 + \epsilon(\epsilon + 2d_{\mathbb{X}}(y, p)) \end{aligned}$$

$d_{\mathbb{X}}(y, p)$ is bounded as P is compact and the relation holds for any ϵ . Thus, f is 1-Lipschitz. ■
When P is a set of points, we build power diagrams with the same construction as Voronoi diagrams, but using power cells instead of Voronoi cells.

Definition 3.27 *Let f be a power distance associated with (P, w) in a metric space \mathbb{X} . For any $p \in P$, the power cell associated with p is defined by:*

$$C(p) = \left\{ x \in \mathbb{X} \mid \forall q \in P, \sqrt{d_{\mathbb{X}}(p, x)^2 + w_p^2} \leq \sqrt{d_{\mathbb{X}}(q, x)^2 + w_q^2} \right\}$$

In Euclidean spaces, any affine diagram can be expressed in the way of a power diagram defined with a power distance. This result is due to Aurenhamer and Imai [7, Theorem 3]. In our case, this translates in the following theorem, which is a restriction of [25, Proposition 3.1]

Theorem 3.28 *Let μ be a probability measure with finite support on an Euclidean space \mathbb{R}^d and $m \in [0, 1[$ be a mass parameter. There exists a set of points P and of weights $(w_p)_{p \in P}$ such that $d_{\mu, m} = f$, where f is the power distance associated with (P, w) .*

The power distance f can be constructed as shown in [67] and detailed in section 4.1. The interesting part is that the sub-level sets of a power distance are a union of balls centred on the points of P . Thus we can describe the sub-level sets of $d_{\mu, m}$ and use the Nerve Theorem to build a simplicial complex similar to the Čech complex to compute the persistent homology of $d_{\mu, m}$. However, the number of balls is the same as the number of non-empty k^{th} -order Voronoi cells in the case of an empirical measure. This number can be of order $O\left(n^{\lfloor \frac{d+1}{2} \rfloor} k^{\lceil \frac{d+1}{2} \rceil}\right)$ [39] which make it impossible to use in practice. Chapters 4 and 5 will tackle this problem using approximating schemes.

4 Distance to a measure approximation

The interest of the distance to a measure for handling noise in persistent homology is hindered by the complexity of its sub-level sets. Theorem 3.28 guarantees that in Euclidean spaces, the sub-level sets can be described as a union of balls. However, the number of balls is too large to make it usable in practice. Moreover, this description cannot be used in other metric spaces. This chapter explores schemes to approximate the distance to a measure in order to reduce the number of balls needed to describe the sub-level set of the distance to a measure.

First, we consider the case of an empirical measure μ on a finite point cloud P in an Euclidean space. The results are provided without any assumption on the nature of μ or the generative model of P . In a second time, we provide a way to approximate the distance to an empirical measure in any metric space with a linear number of balls. Moreover, this approximation scheme can be extended to any probability measure with compact support.

Let us make a quick recall of notation. When working with empirical measures, we consider a point set P with n points and a mass parameter m . $k = mn$ is assumed to be an integer for the sake of simplicity. All results adapt if it is not the case but all proofs need specific considerations for the $\lceil k \rceil^{th}$ -nearest neighbour as shown in the proofs of Propositions 3.13 and 3.14.

4.1 Barycentric decomposition

In Euclidean spaces, the distance to an empirical measure can be described as a power distance following Theorem 3.28. This means that its sub-level sets are a union of balls. The exact expression of the power distance can be easily computed. Its support is a subset of $\Lambda_k(P)$, the set of all barycentres of k points of P . This observation was first made in [67]. Let us write $(p_1(x), \dots, p_k(x))$ the k -nearest neighbours of x and let $bar(x)$ be their barycentre. Then the distance to a measure becomes:

$$\begin{aligned}
d_{\mu,m}(x)^2 &= \frac{1}{k} \sum_{i=1}^k \|p_i(x) - x\|^2 \\
&= \frac{1}{k} \sum_{i=1}^k (\|p_i(x) - \text{bar}(x)\|^2 + \|\text{bar}(x) - x\|^2 + 2 \langle p_i(x) - \text{bar}(x) | \text{bar}(x) - x \rangle) \\
&= \frac{1}{k} \sum_{i=1}^k \|p_i(x) - \text{bar}(x)\|^2 + \|\text{bar}(x) - x\|^2 + 2 \langle \frac{1}{k} \sum_{i=1}^k p_i(x) - \text{bar}(x) | \text{bar}(x) - x \rangle \\
&= \frac{1}{k} \sum_{i=1}^k \|p_i(x) - \text{bar}(x)\|^2 + \|\text{bar}(x) - x\|^2
\end{aligned}$$

Remark that the first expression in the last equality depends only on the k^{th} -order Voronoi cell in which the point x is located. These cells are each associated with a set of k points of P and thus with a barycentre. Let us introduce the weight associated with the barycentre, which is the weight used in the power distance.

Definition 4.1 Let (q_1, \dots, q_k) be a set of k points of P , then the weight associated with $y = \frac{1}{k} \sum_{i=1}^k q_i$ is:

$$w_y^2 = \frac{1}{k} \sum_{i=1}^k \|q_i - y\|^2$$

Proposition 4.2 Let P be a point cloud in \mathbb{R}^d and $m \in [0, 1[$ be a mass parameter such that $k = mn$. The distance to μ , $d_{\mu,m}$, is equal to the power distance associated with $(\Lambda_k(P), w)$, where w is the weight from the previous definition. It means that:

$$d_{\mu,m}(x) = \min_{y \in \Lambda_k(P)} \sqrt{\|x - y\|^2 + w_y^2} = \sqrt{\|\text{bar}(x) - x\|^2 + w_{\text{bar}(x)}^2}.$$

Proof: Let x be a point of \mathbb{R}^d . We have $d_{\mu,m}(x)^2 = w_{\text{bar}(x)}^2 + \|\text{bar}(x) - x\|^2$ and $\text{bar}(x) \in \Lambda_k(P)$.

Thus $d_{\mu,m}(x) \geq \min_{y \in \Lambda_k(P)} \sqrt{\|x - y\|^2 + w_y^2}$.

Now, consider a set of points $(q_1, \dots, q_k) \in P^k$ and y their barycentre.

$$\|x - y\|^2 + w_y^2 = \frac{1}{k} \sum_{i=1}^k \|q_i - x\|^2 \geq \frac{1}{k} \sum_{i=1}^k \|p_i(x) - x\|^2 = d_{\mu,m}(x)^2$$

Hence $d_{\mu,m}(x) = \min_{y \in \Lambda_k(P)} \sqrt{\|x - y\|^2 + w_y^2}$. ■

The sub-level sets of $d_{\mu,m}$ are thus a union of at most $\binom{n}{k}$ balls. In practice, some of these balls are always included in others. They corresponds to sets of k points of P such that no point $x \in \mathbb{R}^d$ has them as its k -nearest neighbours. In other words, their k^{th} -order Voronoi cell is empty. The set of barycentres $\Lambda_k(P)$ can thus be trimmed to the set of barycentres corresponding to non-empty k^{th} -order Voronoi cell without loss of precision. Although this reduction is welcome, it is not sufficient to make computation tractable as the number of non-empty cells can be $O\left(n^{\lfloor \frac{d+1}{2} \rfloor} k^{\lceil \frac{d+1}{2} \rceil}\right)$ [39].

4.2 Random sampling

Reducing the number of balls needed to describe the sub-level sets of the distance to a measure is a necessity to be able to compute persistence diagrams in reasonable time. This can be seen as reducing the number of points supporting the power distance. A natural way to do it is to sample the set of barycentres.

We already saw that some of the barycentres are useless. Thus, we want to sample among the other ones, those that are associated with a non-empty Voronoi cell. The 1-Lipschitz property of $d_{\mu,m}$ implies that removing a barycentre with a small cell will create an error of at most the diameter of the cell.

Lemma 4.3 *Given P a finite point set of \mathbb{R}^d , a fixed $m \in [0, 1[$ and μ the empirical measure on P , let S be the support of $d_{\mu,m}$ and w the weight function. For any $b \in S$, the power distance f associated with $(S \setminus \{b\}, w|_{S \setminus \{b\}})$ satisfies $\|f - d_{\mu,m}\|_\infty \leq \Delta_b$, where Δ_b is the diameter of the Voronoi cell associated with b .*

Proof: Let C be the cell associated with b . Remark that f and $d_{\mu,m}$ are equal outside C . Take $x \in C$. There exists a point $y \notin C$ such that $d_{\mathbb{X}}(x, y) \leq \frac{\Delta_b}{2}$. $d_{\mu,m}$ and f are 1-Lipschitz functions and hence $|f(x) - d_{\mu,m}(x)| \leq 2d_{\mathbb{X}}(x, y) \leq \Delta_b$. ■

Selecting barycentres associated with bigger cells is hence natural. However, if there is an area with only small cells, it is necessary to take at least one of the barycentres. We propose an algorithm that samples randomly the barycentres. Given a point cloud P and a number i of queries, it outputs a set of barycentres and their weights. Each query corresponds to a random choice of a point and the computation of the barycentre associated with the cell containing the point. The algorithm needs a box enclosing the point cloud P . It can be interesting to use a box larger than the minimal enclosing one due to some outer Voronoi cell we may want to capture.

$d_{\mu,m}$ APPROXIMATION BY RANDOM SAMPLING

1. Compute an enclosing box of P .
 2. Repeat i times:
 - (a) Pick uniformly at random a point x inside the box.
 - (b) Compute the barycentre y of the k nearest neighbours of x .
 - (c) If y is not yet in the result set B , add it to it and compute its weight.
-

We first show how the choice of the enclosing box affects the approximation. Let us write $d_{\mu,m}^B$, the approximation obtained by sampling the barycentres.

Lemma 4.4 *If $d_{\mu,m}^B$ approximates $d_{\mu,m}$ with additive precision ϵ in the enclosing box \mathcal{B} , then*

for any $x \in \mathbb{R}^d$,

$$|d_{\mu,m}^B(x) - d_{\mu,m}(x)| \leq \min(\max_{y \in \partial \mathcal{B}} (d_{\mu,m}(y)) + \epsilon; \epsilon + 2d_{\mathbb{X}}(x, \mathcal{B})),$$

where $\partial \mathcal{B}$ is the boundary of \mathcal{B} .

Proof: Trivially, we have $d_{\mu,m}^B \geq d_{\mu,m}$. Consider a point x at distance δ from the enclosing box \mathcal{B} and let $\pi(x)$ be its projection onto \mathcal{B} .

$$\delta \leq d_{\mu,m}(x) \leq d_{\mu,m}^B(x) \leq d_{\mu,m}^B(\pi(x)) + \delta \leq d_{\mu,m}(\pi(x)) + \epsilon + \delta \leq d_{\mu,m}(x) + \epsilon + 2\delta.$$

■

The choice of the enclosing box must balance $\max_{y \in \partial \mathcal{B}} d_{\mu,m}(y)$ that increases as the box grows and $d_{\mathbb{X}}(x, \partial \mathcal{B})$ that decreases. Now, we prove that we are able to approximate $d_{\mu,m}$ with precision ϵ inside the enclosing box. We recall that $V_d^0(r)$ is the volume of a ball of radius r in \mathbb{R}^d .

Theorem 4.5 *Let P be a finite point set in \mathbb{R}^d and $m \in [0, 1[$ be a mass parameter. Considering the empirical measure μ on P , the random sampling algorithm returns an ϵ additive approximation of $d_{\mu,m}$ in the enclosing box \mathcal{B} after $O\left(n^{\lfloor \frac{d+1}{2} \rfloor} k^{\lceil \frac{d+1}{2} \rceil} + d \log(n+k) \frac{\Delta}{V_d^0(\epsilon)}\right)$ queries with probability at least $\frac{1}{2}$, where Δ is the volume of \mathcal{B} .*

Proof: Remark that given a k^{th} -order Voronoi cell V , if one of the query point was in V then $d_{\mu,m}|_V = d_{\mu,m}^B|_V$. Let T be the union of all such cells. Then for all $x \in \mathbb{R}^d$ such that $B(x, \frac{\epsilon}{2}) \cap T \neq \emptyset$, $|d_{\mu,m}(x) - d_{\mu,m}^B(x)| \leq \epsilon$ because $d_{\mu,m}$ and $d_{\mu,m}^B$ are 1-Lipschitz.

We bound the number of queries needed such that for all $x \in \mathcal{B}$, $B(x, \frac{\epsilon}{2}) \cap T \neq \emptyset$. We introduce $p = \left(\frac{V_d^0(\frac{\epsilon}{2})}{\Delta}\right)$, the probability for a query point to be in a given ball of radius $\frac{\epsilon}{2}$.

We consider an uniform sampling of i queries, X_1, \dots, X_i . We now build a tree where each node corresponds to an event with some probability. We start with $T_1 = \emptyset$. Assuming that there exists a x_1 such that $B(x_1, \frac{\epsilon}{2}) \cap T_1 = \emptyset$, we denote as A_1 the event, $\exists j \leq i, X_j \in B(x_1, \frac{\epsilon}{2})$. Remark that $Pr(A_1) \geq 1 - (1-p)^i$. The ball $B(x_1, \frac{\epsilon}{2})$ intersects one or more k^{th} -order Voronoi cells. We denote by $C_1^1, \dots, C_1^{c_1}$ these cells. Let B_1^1 the event that the first point to hit $B(x_1, \frac{\epsilon}{2})$ is in C_1^1 . In this case, $T_2 = T_1 \cup C_1^1$. If there exists x_2 such that $B(x_2, \frac{\epsilon}{2}) \cap T_2 = \emptyset$, then we start again with x_2 . The probability that one query point is inside $B(x_2, \frac{\epsilon}{2})$ conditionally to B_1^1 is greater than $1 - (1-p)^{i-|T_2|}$.

Recursively, we can define the whole tree. The depth of the tree is bounded by the number of the number \mathcal{N} of non-empty k^{th} -order Voronoi cells. Given a path from the root to a leaf, the realisation of all events on this path guarantees that the propriety needed for the theorem is realised.

At each node, the conditional probability of the event A conditionally to the event B is bounded from below by $1 - (1-p)^{i-|\mathcal{N}|}$. Moreover, the union of the event of type B conditionally to the event of type A is 1. By recursion, we can go back from the leaves to the root. The probability that at least one of the path has all events realised is at least $1 - \mathcal{N}(1-p)^{i-|\mathcal{N}|}$. Knowing

that $\mathcal{N} = O\left(n^{\lfloor \frac{d+1}{2} \rfloor} k^{\lceil \frac{d+1}{2} \rceil}\right)$, and to guarantee that this probability is more than $\frac{1}{2}$, we need $i = O\left(n^{\lfloor \frac{d+1}{2} \rfloor} k^{\lceil \frac{d+1}{2} \rceil} + d \log(n+k) \frac{\Delta}{V_d^0(\epsilon)}\right)$.

Remark that the condition $B(x, \frac{\epsilon}{2}) \cap T \neq \emptyset$ is necessary to ensure that the error in x is less than ϵ . Consider the power distance f associated with $(Q, 0)$, where Q is composed of x and an arbitrary number of points located on the sphere of centre x and radius ϵ . This setting can be achieved as a distance to a measure by considering the empirical measure on Q and $k = 1$. Then, if there is no query point inside $B(x, \frac{\epsilon}{2})$, x is not a selected barycentre and $d_{\mu, m}^B(x) = \epsilon$, while $d_{\mu, m}(x) = 0$. ■

Although, this approximation scheme is working rather well in practice with a number of queries linear in the size of P , we are unable to give guarantees on the result due to the randomness. Remark that the number of queries needed to obtain guarantees is of the same order as the number of points needed by the algorithm taking the queries on a grid.

We now concentrate on linear sized sets to approximate the distance to a measure using deterministic algorithms. These schemes provide multiplicative approximation bounds instead of additive ones.

4.3 Witnessed k -distance

The first method to approximate the distance $d_{\mu, m}$ to an empirical measure μ in an Euclidean space was proposed in [67]. Using the power distance structure of $d_{\mu, m}$, we select only the barycentres that are associated with a cell containing a point of P . On the first hand, this guarantees that we select at most $|P|$ points to describe the new power distance. On the other hand, it makes sense for inference purpose because P is supposed to sample an underlying object. We want to recover precisely $d_{\mu, m}$ in the area located near the underlying object and, hence, near P .

Definition 4.6 *Let P be a finite point set of \mathbb{R}^d and let $m \in [0, 1[$ be a mass parameter. Considering the empirical measure μ over P , the witnessed k -distance is defined as*

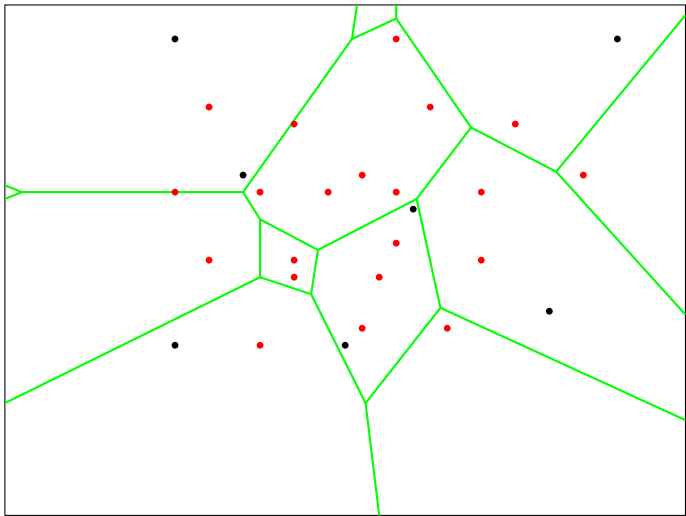
$$d_{\mu, m}^W(x) = \min_{p \in P} \sqrt{w_{bar(p)}^2 + ||bar(p) - x||^2}.$$

The computation of the weights does not present more difficulty than the search of the k nearest neighbours for each point. The number of queries is exactly the size of the point set P . Figure 4.1 shows a point set P and the barycentres for $k = 2$. From top to bottom, it displays the set of all barycentres, then the set of all barycentres associated with a non-empty cell and then the set of barycentres kept when using the witnessed k -distance.

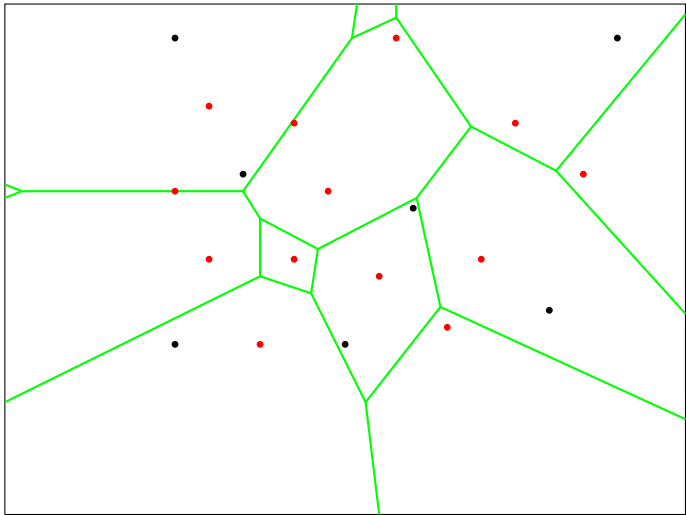
A bound on the quality of the approximation was given [67, Lemma 3.3]. We improve this bound to be at least as good as our new approximation, described in Section 4.4. We are not able to prove the tightness of the bound. However, we can give a lower bound on the precision.

Theorem 4.7 *Let P be a finite point set of \mathbb{R}^d and let $m \in [0, 1[$ be a mass parameter. Considering the empirical measure μ over P ,*

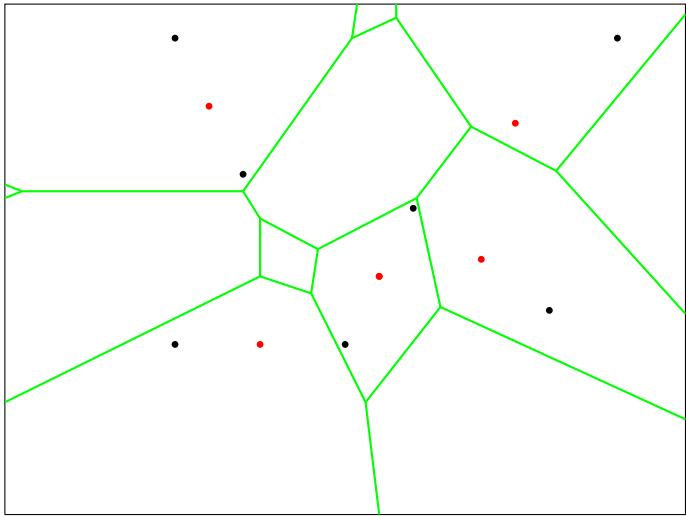
$$d_{\mu, m} \leq d_{\mu, m}^W \leq \sqrt{6} d_{\mu, m}.$$



(a) All barycentres for $k = 2$



(b) Barycentres associated with non-empty cells for $k = 2$



(c) Barycentres used by the witnessed k -distance

Figure 4.1 – Set of barycentres associated with a point cloud

The previous version of this theorem [67, Lemma 3.3] used a 3 instead of the $\sqrt{6}$. We postpone the proof as it will be a quite direct corollary of Theorem 4.15 and concentrate first on studying the tightness of the bound.

The tightness of the lower bound is obvious as it suffices to take $k = 1$ to get an equality between $d_{\mu,m}$ and $d_{\mu,m}^W$. However, we do not know if the upper bound is tight. The bound $\sqrt{6}$ can not be less than $1 + \sqrt{2}$, whose value is greater than $\sqrt{5.82}$.

Let us introduce the following example in \mathbb{R}^d . We fix $k = 2d$ and $0 < \epsilon < \sqrt{2}$. The point cloud P consists of $4d^2$ points located at the coordinates $(0, \dots, 0, \alpha, 0, \dots, 0)$ with multiplicity 1 when $\alpha = 1$ or $\alpha = -1$ and multiplicity $2d - 1$ when $\alpha = 1 + \sqrt{2} - \epsilon$ or $\alpha = \epsilon - 1 - \sqrt{2}$. Figure 4.2 gives its representation in dimension 2 where triangles have multiplicity 1 and circles have multiplicity 3.

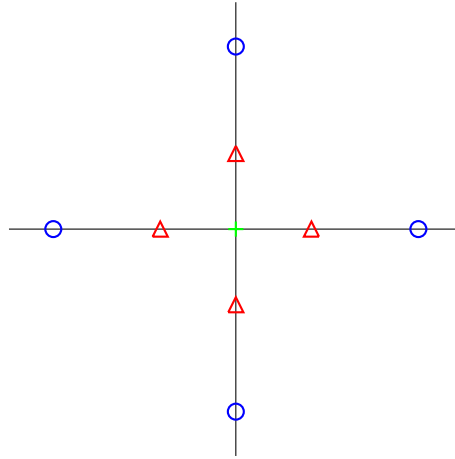


Figure 4.2 – Example for the tightness of $d_{\mu,m}^W$

The points are placed such that the k nearest neighbours of any triangle are itself and the $k - 1$ points located at the nearest circle. These k nearest neighbours are also the ones from the circles.

Consider the value of the $d_{\mu,m}$ and $d_{\mu,m}^W$ at the origin o . Each of the k nearest neighbors of o is at distance exactly 1 from o . Hence,

$$d_{\mu,m}(o) = 1.$$

The construction induced that the structure is perfectly symmetric and the set of barycentres W we consider in the witnessed k -distance contains exactly $2d$ points. These points are located at the coordinates $(0, \dots, 0, \alpha, 0, \dots, 0)$ where $\alpha = 1 + \frac{2d-1}{2d}(\sqrt{2} - \epsilon)$ or the opposite.

Let b be a member of W . The weight associated with b is

$$\begin{aligned} w_b &= \frac{1}{2d} \left[\left(\frac{2d-1}{2d} (\sqrt{2}-\epsilon) \right)^2 + (2d-1) \left(\frac{1}{2d} (\sqrt{2}-\epsilon) \right)^2 \right] \\ &= \frac{2d-1}{(2d)^3} \left[(2d-1)(\sqrt{2}-\epsilon)^2 + (\sqrt{2}-\epsilon)^2 \right] \\ &= \frac{2d-1}{(2d)^2} (\sqrt{2}-\epsilon)^2. \end{aligned}$$

All of the points of W are located at the same distance to o . Hence,

$$\begin{aligned} d_{\mu,m}^W(o)^2 &= w_b + \left(1 + \frac{2d-1}{2d} (\sqrt{2}-\epsilon) \right)^2 \\ &= \frac{2d-1}{(2d)^2} (\sqrt{2}-\epsilon)^2 + 1 + \frac{2d-1}{d} (\sqrt{2}-\epsilon) + \frac{(2d-1)^2}{(2d)^2} (\sqrt{2}-\epsilon)^2 \\ &= \frac{1}{2d} + \frac{2d-1}{2d} \left(1 + 2(\sqrt{2}-\epsilon) + (\sqrt{2}-\epsilon)^2 \right) \\ &= \frac{1}{2d} + \frac{2d-1}{2d} (1 + \sqrt{2}-\epsilon)^2. \end{aligned}$$

Since we can take ϵ as small as we want and make the dimension grow, this relation ensures that we cannot find a better constant than $1 + \sqrt{2}$ in Theorem 4.7. Remark that this does not change if we apply a homothety to the example. Hence an additive bound is out of reach. Previous work in [67] concentrated on the inference purpose of the distance to a measure, which is not our focus here. Assuming the presence of a ground truth under the guise of a Riemannian manifold, it is possible to achieve some additive bound between the distance to the manifold and the witnessed k -distance.

4.4 Power distance with compact support

The witnessed k -distance presents two properties that we want to get rid of. First, it is only usable in Euclidean spaces as it involves barycentres. For more general classes of spaces, barycentres do not always exist and the distance to an empirical measure can not necessarily be expressed as a power distance. The second point is related to scalar field analysis. Assume that the points of P are sampled on a manifold M and that a real valued function f is defined on the manifold. If we are interested in analysing the structure of f , it does not make a lot of sense to consider the barycentres that can be outside M and it is difficult to assign them a value.

The functions used until now are power distances f associated with (Q, w) , where Q is a subset of $\Lambda_k(P)$, the set of all barycentres of k points of P . Therefore the sub-level set $f^{-1}(]-\infty, \alpha])$ is a union of balls centred on the points q of Q . Their radius can be computed and is given by $r_q(\alpha) = \sqrt{\alpha^2 - w_q^2}$. When this radius is imaginary, the ball is considered empty by convention. Power distances are stable under small perturbations of the points and thus we can hope to change the support to obtain a more interesting approximation of $d_{\mu,m}$.

We first discuss the stability of power distances and then provide an approximation to $d_{\mu,m}$,

whose support is exactly P and give guarantees equivalent to the one obtained with the the witnessed k -distance.

4.4.1 Stability of power distances

The stability of power distances is directly linked to the weights put on the points. Considering two weighted points clouds P and Q in a metric space \mathbb{X} , we assume that the weights of the power distances are t -Lipschitz. This means that they are induced by a t -Lipschitz function on \mathbb{X} and are obtained by restricting this function to P and Q . Stability is then obtained with respect to the Hausdorff distance between P and Q .

We present two different results. Proposition 4.10 is better than Proposition 4.8 but maybe less straightforward.

Proposition 4.8 *Let \mathbb{X} be a metric space, and let $w : \mathbb{X} \rightarrow \mathbb{R}$ be a function from \mathbb{X} to \mathbb{R} . Let P and Q be two compact subsets of \mathbb{X} . Let f_P and f_Q be the power distances associated with $(P, w|_P)$ and $(Q, w|_Q)$. If w is t -Lipschitz with $t \geq 1$, then:*

$$\|f_P - f_Q\|_\infty \leq \sqrt{2} t d_H(P, Q).$$

Proof: Let x be a point of \mathbb{X} and $q \in Q$ such that $f_Q(x)^2 = w_q^2 + d_{\mathbb{X}}(x, q)^2$. There exists a point $p \in P$ such that $d_{\mathbb{X}}(p, q) = \epsilon \leq d_H(P, Q)$. So,

$$\begin{aligned} f_P(x)^2 &\leq w_p^2 + d_{\mathbb{X}}(p, x)^2 \\ &\leq (w_q + t\epsilon)^2 + (d_{\mathbb{X}}(q, x) + \epsilon)^2 \\ &= w_q^2 + d_{\mathbb{X}}(q, x)^2 + 2(tw_q + d_{\mathbb{X}}(q, x))\epsilon + (1 + t^2)\epsilon^2. \end{aligned}$$

Using $t \geq 1$,

$$f_P(x)^2 \leq f_Q(x)^2 + 2t(w_Q + d_{\mathbb{X}}(q, x))\epsilon + 2t^2\epsilon^2.$$

Moreover, the relation $a + b \leq \sqrt{2}\sqrt{a^2 + b^2}$ implies

$$f_P(x)^2 \leq f_Q(x)^2 + 2\sqrt{2}t f_Q(x)\epsilon + 2t^2\epsilon^2 = (f_Q(x) + \sqrt{2}t\epsilon)^2.$$

To conclude the proof, it suffices to reverse the roles of P and Q . ■

The second stability result relies on a lemma about inclusions between balls. In addition to the stability result of Proposition 4.10, it is also useful for the study of the weighted Rips filtration of Chapter 5.

Lemma 4.9 *Let $p, q \in \mathbb{X}$ be points such that $d_{\mathbb{X}}(p, q) \leq \epsilon$, and let $w : \mathbb{X} \rightarrow \mathbb{R}$ be a t -Lipschitz function. For all $\alpha \geq w_p$,*

$$r_p(\alpha) + \epsilon \leq r_q(\alpha + \sqrt{1 + t^2} \epsilon).$$

Proof: First, observe that $r_p(\alpha)$ can be bounded by

$$\begin{aligned} r_p(\alpha)^2 &= \alpha^2 - w_p^2 \leq \alpha^2 - w_p^2 + (t\alpha - \sqrt{1+t^2} w_p)^2 \\ &= (\sqrt{1+t^2} \alpha - t w_p)^2. \end{aligned}$$

Next, we relate r_p and r_q ,

$$\begin{aligned} (r_p(\alpha) + \epsilon)^2 &= \alpha^2 - w_p^2 + 2\epsilon\sqrt{\alpha^2 - w_p^2} + \epsilon^2 \\ &\leq \alpha^2 - w_p^2 + 2\epsilon(\sqrt{1+t^2} \alpha - t w_p) + \epsilon^2 \\ &= (\alpha + \sqrt{1+t^2} \epsilon)^2 - (w_p + t\epsilon)^2 \\ &\leq (\alpha + \sqrt{1+t^2} \epsilon)^2 - w_q^2 \\ &= r_q(\alpha + \sqrt{1+t^2} \epsilon)^2. \end{aligned}$$

The requirement that $\alpha \geq w_p$ allows us to take the square root of both sides of the inequality since both will be nonnegative. ■

The lemma means that for a parameter α , the ball centred at p with radius $r_p(\alpha)$ is included inside a ball centred at q with radius α' , slightly larger than α , as illustrated by Figure 4.3.

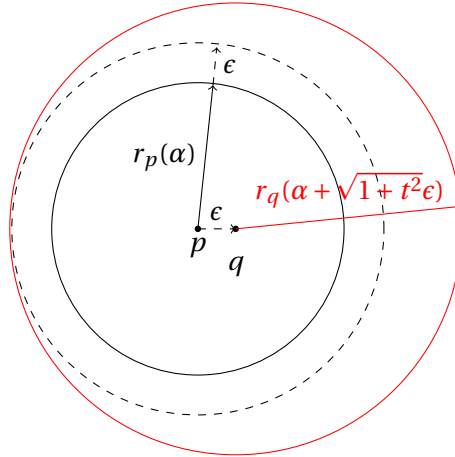


Figure 4.3 – Inclusion of weighted balls

Proposition 4.10 *Let \mathbb{X} be a metric space and let $w : \mathbb{X} \rightarrow \mathbb{R}$ be a function. Let P and Q be two compact subsets of \mathbb{X} . Let f_P and f_Q be the power distances associated with $(P, w|_P)$ and $(Q, w|_Q)$. If w is t -Lipschitz, then*

$$\|f_P - f_Q\|_\infty \leq \sqrt{1+t^2} d_H(P, Q).$$

Proof: Let x be any point of \mathbb{X} . There exists $p \in P$ such that $x \in \bar{B}(p, r_p(f_P(x)))$. There also exists $q \in Q$ such that $d_{\mathbb{X}}(p, q) \leq d_H(P, Q)$. By Lemma 4.9 and the triangle inequality, $x \in \bar{B}(q, r_q(f_P(x) + \sqrt{1+t^2} d_H(P, Q)))$. Thus, $f_Q(x) \leq f_P(x) + \sqrt{1+t^2} d_H(P, Q)$. P and Q are interchangeable, therefore $\|f_Q - f_P\|_\infty \leq \sqrt{1+t^2} d_H(P, Q)$. ■

4.4.2 Approximation of the distance to a measure

Given a probability measure μ defined over a metric space \mathbb{X} , we propose an approximation $d_{\mu,m}^P$ to $d_{\mu,m}$ using power distances. Remark that we are no longer restricted to empirical measures. This approximation starts by fixing a set P , which will serve as the support of the power distance. We show that, for a well chosen P , we obtain a good approximation of $d_{\mu,m}$.

Definition 4.11 *Let μ be a probability measure on a metric space \mathbb{X} and let $m \in [0, 1[$ be a mass parameter. Given a subset P of \mathbb{X} , we define $d_{\mu,m}^P$ as the power distance associated with $(P, d_{\mu,m})$:*

$$d_{\mu,m}^P(x) = \sqrt{\min_{p \in P} d_{\mu,m}(p)^2 + d_{\mathbb{X}}(p, x)^2}.$$

The weight of each point is its distance to the measure μ . If P is close to $\text{Supp}(\mu)$, we obtain an approximation of $d_{\mu,m}$.

Theorem 4.12 *Let μ be a probability measure on a metric space \mathbb{X} and let $m \in [0, 1[$ be a mass parameter. Let P be a subset of \mathbb{X} . If P is an ϵ -sample of $\text{Supp}(\mu)$, then*

$$\frac{1}{\sqrt{2}} d_{\mu,m} \leq d_{\mu,m}^P \leq \sqrt{5} (d_{\mu,m} + \epsilon).$$

Proof: Let x be a point of \mathbb{X} . Using notations of Section 3.4,

$$d_{\mu,m}(x)^2 = \frac{1}{m} \int_{\mathbb{X}} d_{\mathbb{X}}(y, x)^2 \mu_{x,m}(y) dy.$$

Fix a point $p \in P$. Since $\mu_{p,m}$ is a submeasure of μ of total mass m ,

$$\begin{aligned} d_{\mu,m}(x)^2 &= \frac{1}{m} \int_{\mathbb{X}} d_{\mathbb{X}}(y, x)^2 \mu_{x,m}(y) dy \\ &\leq \frac{1}{m} \int_{\mathbb{X}} d_{\mathbb{X}}(y, x)^2 \mu_{p,m}(y) dy \\ &\leq \frac{1}{m} \int_{\mathbb{X}} (d_{\mathbb{X}}(y, p) + d_{\mathbb{X}}(p, x))^2 \mu_{p,m}(y) dy \\ &\leq d_{\mathbb{X}}(p, x)^2 \frac{2}{m} \int_{\mathbb{X}} \mu_{p,m}(y) dy + \frac{2}{m} \int_{\mathbb{X}} d_{\mathbb{X}}(y, p)^2 \mu_{p,m}(y) dy \\ &= 2(d_{\mathbb{X}}(p, x)^2 + d_{\mu,m}(p)^2). \end{aligned}$$

The third inequality follows from the triangle inequality and the relation $(a + b)^2 \leq 2(a^2 + b^2)$. As the inequality holds for any point p in P we conclude that

$$d_{\mu,m}(x) \leq \sqrt{2} d_{\mu,m}^P(x).$$

To show the other inequality, we again fix $p \in P$. By definition,

$$\begin{aligned} d_{\mu,m}^P(x)^2 &\leq d_{\mathbb{X}}(x, p)^2 + d_{\mu,m}(p)^2 \\ &\leq d_{\mathbb{X}}(x, p)^2 + \frac{1}{m} \int_{\mathbb{X}} d_{\mathbb{X}}(p, y)^2 \mu_{x,m}(y) dy \\ &\leq d_{\mathbb{X}}(x, p)^2 + \frac{1}{m} \int_{\mathbb{X}} (d_{\mathbb{X}}(p, x) + d_{\mathbb{X}}(x, y))^2 \mu_{x,m}(y) dy \\ &\leq 3 d_{\mathbb{X}}(x, p)^2 + 2 d_{\mu,m}(x)^2. \end{aligned}$$

By definition of $d_{\mu,m}$, $d_{\mathbb{X}}(x, \text{Supp}(\mu)) \leq d_{\mu,m}(x)$. Consequently, there exists a point $p \in P$ such that $d_{\mathbb{X}}(x, p) \leq d_{\mu,m}(x) + \epsilon$. Hence,

$$d_{\mu,m}^P(x)^2 \leq 5(d_{\mu,m}(x) + \epsilon)^2.$$

■

4.4.3 Restriction to measures with finite support

Given a finite set of points P in a metric space \mathbb{X} , we want to approximate the distance to the empirical measure μ on P . Taking P itself as the support for the power distance, Theorem 4.12 yields an immediate corollary as $\epsilon = 0$. Moreover, these bounds are tight.

Corollary 4.13 *Let P be a finite point set of a metric space \mathbb{X} and $m \in [0, 1[$ be a mass parameter. Considering the empirical measure μ on P ,*

$$\frac{1}{\sqrt{2}} d_{\mu,m} \leq d_{\mu,m}^P \leq \sqrt{5} d_{\mu,m}.$$

Proposition 4.14 *The bounds of Corollary 4.13 are tight.*

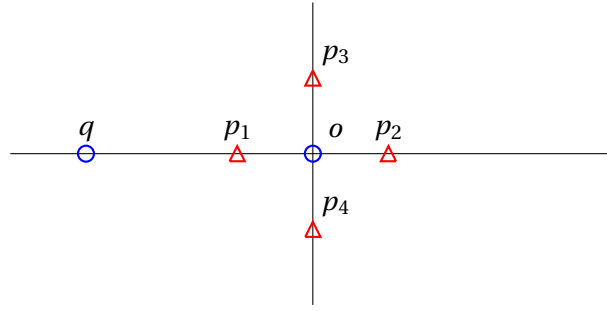
Proof: We are looking for a worst case scenario where inequalities become equalities for at least one point. We consider the space \mathbb{R}^d with the L_1 -norm, denoted $|\cdot|$. For any fixed dimension d , we build the set of $2d$ points whose coordinates have the form $(0, \dots, 0, \pm 1, 0, \dots, 0)$. These points are marked by triangles in the 2-dimensional illustration of Figure 4.4. Remark that their pairwise distances are all equal to 2. We fix $k = 2d$ and we study $d_{\mu,m}$ and $d_{\mu,m}^P$ at points $q(-3, 0, \dots, 0)$ and o . First we compute the value of $d_{\mu,m}(p)$ for any $p \in P$:

$$d_{\mu,m}(p)^2 = \frac{1}{2d} \sum_{q \in P} |q - p|^2 = \frac{1}{2d} \sum_{q \in P \setminus \{p\}} 2^2 = 4 \frac{2d-1}{2d} = 4 - \frac{2}{d}$$

Now we compute the value of $d_{\mu,m}$ at q and o :

$$d_{\mu,m}(o)^2 = \frac{1}{2d} \sum_{p \in P} |p - o|^2 = 1$$

$$d_{\mu,m}(q)^2 = \frac{1}{2d} \sum_{p \in P} |p - q|^2 = \frac{1}{2d} (4 + (2d-1)16) = 16 - \frac{6}{d}$$


 Figure 4.4 – Example for the tightness of $d_{\mu,m}^P$ in a general metric space

All the points p have the same value for $d_{\mu,m}$. Hence,

$$d_{\mu,m}^P(o)^2 = d_{\mu,m}(p)^2 + |p - o|^2 = 5 - \frac{2}{d}$$

$$d_{\mu,m}^P(q)^2 = d_{\mu,m}(p)^2 + |p - q|^2 = 8 - \frac{2}{d}$$

When d increases, the ratio $\frac{d_{\mu,m}^P(o)}{d_{\mu,m}(o)}$ tends to $\sqrt{5}$, while $\frac{d_{\mu,m}^P(q)}{d_{\mu,m}(q)}$ tends to $\frac{1}{\sqrt{2}}$. Thus, the bounds of Corollary 4.13 are reached at the limit for the same data set, although at two different points. Remark that the construction can be arbitrarily scaled thus an additive bound is out of reach. ■

4.4.4 Euclidean case

Restricting the class of metric spaces we consider yields better bounds. We consider the Euclidean space \mathbb{R}^d with the L_2 -norm. Considering a finite point set P and its empirical measure μ in \mathbb{R}^d , we are able to obtain bounds of the same quality as those of the witnessed k -distance.

Theorem 4.15 *Let P be a finite point set in \mathbb{R}^d and let $m \in [0, 1[$ be a mass parameter. Considering the empirical measure μ on P , the following relation is tight.*

$$\frac{1}{\sqrt{2}} d_{\mu,m} \leq d_{\mu,m}^P \leq \sqrt{3} d_{\mu,m}.$$

Proof: The first inequality is exactly the same as the one from Theorem 4.12. For the second inequality, let x be a point in \mathbb{R}^d , and let p be a point of P . Thus,

$$d_{\mu,m}^P(x)^2 \leq d_{\mu,m}(p)^2 + \|p - x\|^2.$$

Using Proposition 4.2,

$$d_{\mu,m}^P(x)^2 \leq w_{bar(x)}^2 + \|p - bar(x)\|^2 + \|p - x\|^2,$$

and with the inner product, this becomes

$$\begin{aligned} d_{\mu,m}^P(x)^2 &\leq w_{bar(x)}^2 + \|x - bar(x)\|^2 + 2\|p - x\|^2 + 2\langle x - bar(x) | p - x \rangle \\ &= d_{\mu,m}(x)^2 + 2\|p - x\|^2 + 2\langle x - bar(x) | p - x \rangle. \end{aligned}$$

Note that

$$2\langle bar(x) - x | x - p \rangle = \|bar(x) - p\|^2 - \|bar(x) - x\|^2 - \|x - p\|^2.$$

Thus,

$$d_{\mu,m}^P(x)^2 \leq d_{\mu,m}(x)^2 + \|p - x\|^2 + \|bar(x) - p\|^2 - \|x - bar(x)\|^2.$$

This relation holds for any point of P . In particular it holds for any of the k nearest neighbours of x . If we take the average over the k nearest neighbours of x and eliminate the negative term $-\|x - bar(x)\|^2$, we obtain

$$d_{\mu_P,m}^P(x)^2 \leq d_{\mu_P,m}(x)^2 + \frac{1}{k} \sum_{p \in NN_k^P(x)} \|p - x\|^2 + \frac{1}{k} \sum_{p \in NN_k^P(x)} \|bar(x) - p\|^2.$$

Using the definitions of weights and of the distance to the measure,

$$d_{\mu,m}^P(x)^2 \leq d_{\mu,m}(x)^2 + d_{\mu,m}(x)^2 + w_{bar(x)}^2$$

where $w_{bar(x)} \leq d_{\mu,m}(x)$. We conclude that

$$d_{\mu,m}^P(x) \leq \sqrt{3} d_{\mu,m}(x).$$

We now provide an example where these bounds are tight. For P , consider the two points a and b on the real line with coordinates 1 and -1 as given in Figure 4.5.

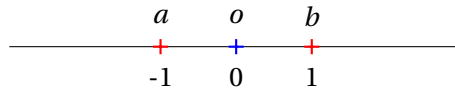


Figure 4.5 – Example for the tightness of $d_{\mu,m}^P$ in Euclidean space

Fix the mass parameter m equal to 1 so that $k = 2$. It follows that

$$d_{\mu,m}(a) = d_{\mu,m}(b) = \sqrt{\frac{1}{2} \|b - a\|^2} = \sqrt{2},$$

$$d_{\mu,m}(o) = \sqrt{\frac{1}{2} \|o - b\|^2 + \|o - a\|^2} = 1.$$

We now compute the last interesting value:

$$d_{\mu,m}^P(o)^2 = d_{\mu,m}(a)^2 + \|a - o\|^2 = 3.$$

We can thus conclude that $d_{\mu,m}^P(o) = \sqrt{3} d_{\mu,m}(o)$ which gives us the tightness of the upper bound.

For the lower, let us modify the weights. a has now weight $(1 - \epsilon)$ and b has weight ϵ . There is only one barycentre c at coordinate $2\epsilon - 1$. Consider the point d at coordinate $3 - 2\epsilon$. Then:

$$d_{\mu,m}(d) = \sqrt{(1 - \epsilon)(4 - 2\epsilon)^2 + \epsilon(2 - 2\epsilon)^2} = \sqrt{16 + 12\epsilon^2}$$

$$d_{\mu,m}^P(d) = \sqrt{(2 - 2\epsilon)^2 + 4(1 - \epsilon)} = \sqrt{8 - 12\epsilon + 4\epsilon^2}$$

ϵ can be taken arbitrarily small and hence the lower bound is tight. ■

Given Theorem 4.15, the proof of Theorem 4.7 becomes easy.

Proof: First remark that $d_{\mu,m}^W$ is a minimum over a smaller set than $d_{\mu,m}$ due to Proposition prop:barDecompo. We thus get $d_{\mu,m} \leq d_{\mu,m}^W$.

Let x be a point in \mathbb{R}^d . Thus for any $p \in P$,

$$\begin{aligned} d_{\mu,m}^W(x)^2 &\leq w_{bar(p)}^2 + ||bar(p) - x||^2 \\ &\leq w_{bar(p)}^2 + ||bar(p) - p||^2 + ||p - x||^2 + 2\langle bar(p) - p | p - x \rangle \\ &\leq d_{\mu,m}(p)^2 + 2||p - x||^2 + ||bar(p) - p||^2 \\ &\leq 2(d_{\mu,m}(p)^2 + ||p - x||^2) \\ &\leq 2 d_{\mu,m}^P(x)^2. \end{aligned}$$

Using Theorem 4.15, we conclude:

$$d_{\mu,m}^W(x) \leq \sqrt{2} d_{\mu,m}^P(x) \leq \sqrt{6} d_{\mu,m}(x). \quad \text{■}$$

4.5 Application to persistence diagrams approximation

We now have a way to efficiently approximate the sub-level sets of the distance to a measure. Thus, we hope to use it to approximate its persistence diagram. The presence of an additive approximation gives a direct way to approximate persistence diagrams.

Lemma 4.16 *Let f and g be two real valued functions with q -tame sub-level sets filtrations on a metric space \mathbb{X} . Then,*

$$d_B(\text{Dgm}(f), \text{Dgm}(g)) \leq ||f - g||_\infty$$

Proof: Recall that $\text{Dgm}(f)$ is the persistence diagram of the sub-level sets filtration $\{f^{-1}([-\infty, \alpha])\}_{\alpha \in \mathbb{R}}$. For any $\alpha \in \mathbb{R}$,

$$f^{-1}([-\infty, \alpha]) \subset g^{-1}([-\infty, \alpha + ||f - g||_\infty]).$$

As the relation is symmetric, this means that the two sub-level sets filtrations are $||f - g||_\infty$ -interleaved and Corollary 2.27 applies. ■

Here, we are interested in a restricted class of functions. We consider the distance to a probability measure μ and some power distances with finite support. First, we need to prove that they have well-defined persistence diagrams. We show that the functions are q -tame using the following theorem:

Theorem 4.17 ([26, Theorem 2.22]) *Let C be a finite simplicial complex and let $f : C \rightarrow \mathbb{R}$ be a continuous function. Then the sub-level sets filtration of f is q -tame.*

First, we show that the distance to a probability measure is q -tame. Remark that $d_{\mu,m}$ is a non-negative function. Thus, we consider the sub-level sets $d_{\mu,m}^{-1}([0, \alpha])$.

Proposition 4.18 *Let μ be a probability measure on a triangulable metric space \mathbb{X} , and let $m \in [0, 1[$ be a mass parameter. The sub-level sets filtration of $d_{\mu,m}$ is q -tame.*

Proof: Following Corollary 3.16, $d_{\mu,m}$ is continuous and non-negative. If for any α , the sub-level set $d_{\mu,m}^{-1}([0, \alpha])$ is compact, then the sub-level sets filtration of $d_{\mu,m}$ is q -tame. Since \mathbb{X} is triangulable, there exists a homeomorphism h from \mathbb{X} to a locally finite simplicial complex C . Then for any $\alpha > 0$, we can restrict the simplicial complex C to a finite simplicial complex C_α that contains the compact $h(d_{\mu,m}^{-1}([0, \alpha]))$. The function $d_{\mu,m} \circ h^{-1}|_{C_\alpha}$ is continuous on C_α . Thus its sub-level sets filtration is q -tame [26, Theorem 2.22].

The construction extends to any α and therefore the sub-level sets filtration of $d_{\mu,m} \circ h^{-1}$ is q -tame. Furthermore, homology is preserved by homeomorphisms and thus we can say that the sub-level sets filtration of $d_{\mu,m}$ is q -tame.

We only need to show that any sub-level set $d_{\mu,m}^{-1}([0, \alpha])$ is compact. Suppose for contradiction that there exists an $\alpha > 0$ such that $d_{\mu,m}^{-1}([0, \alpha])$ is not compact. Then there exists a sequence $(x_i)_{i>0}$ of points of $d_{\mu,m}^{-1}([0, \alpha])$ such that $d_{\mathbb{X}}(x_0, x_n) \rightarrow \infty$ when $n \rightarrow \infty$. Hence we can extract a sub-sequence $(x_{\phi(i)})_{i>0}$ such that, for any i and j , $\bar{B}(x_{\phi(i)}, \sqrt{2}\alpha) \cap \bar{B}(x_{\phi(j)}, \sqrt{2}\alpha) = \emptyset$.

$$d_{\mu,m}(x_{\phi(i)})^2 = \frac{1}{m} \int_0^m \delta_{\mu,l}(x_{\phi(i)})^2 dl \leq \alpha^2.$$

The function $\delta_{\mu,l}(x_{\phi(i)})$ is non-negative and increasing with l and therefore $\frac{m}{2} \delta_{\mu, \frac{m}{2}}(x_{\phi(i)})^2 \leq m\alpha^2$. Using the definition of $\delta_{\mu,m}$, this implies that $\mu(\bar{B}(x_{\phi(i)}, \sqrt{2}\alpha)) \geq \frac{m}{2}$. Measures are σ -additive, so

$$\mu(\mathbb{X}) \geq \sum_{i>0} \mu(\bar{B}(x_{\phi(i)}, \sqrt{2}\alpha)) \geq \sum_{i>0} \frac{m}{2} = \infty.$$

However, μ is a probability measure and therefore $\mu(\mathbb{X}) = 1$. This contradiction implies that $d_{\mu,m}^{-1}([0, \alpha])$ is compact. ■

The power distances with a finite support have the same tameness property.

Proposition 4.19 *Let \mathbb{X} be a triangulable metric space and f be a power distance associated with (P, w) , where P is a finite point set. Then the sub-level sets filtration of f is q -tame.*

Proof: By construction, f is non-negative and continuous and, for any α , the sub-level set $f^{-1}([0, \alpha])$ is a finite union of balls. The space \mathbb{X} is triangulable and hence the sub-level sets of

f are compact. The same construction as the one used for $d_{\mu,m}$ works and we obtain that f is q -tame. ■

In our approximation, we obtain a multiplicative approximation of $d_{\mu,m}$ and not an additive one. This kind of approximation gives an interleaving on a logarithmic scale.

Proposition 4.20 *Let P be a finite point set of a triangulable metric space \mathbb{X} and $m \in [0, 1[$ be a mass parameter. Considering the empirical measure μ on P ,*

$$d_B^{\log}(\text{Dgm}(d_{\mu,m}), \text{Dgm}(d_{\mu,m}^P)) \leq \ln(\sqrt{5}).$$

Proof: Corollary 4.13 implies that

$$\ln(d_{\mu,m}) - \ln(\sqrt{2}) \leq \ln(d_{\mu,m}^P) \leq \ln(\sqrt{5}) + \ln(d_{\mu,m}).$$

The sub-level sets of $\ln(d_{\mu,m})$ and $\ln(d_{\mu,m}^P)$ are thus $\ln(\sqrt{5})$ -interleaved and Theorem 2.23 applies. ■

Proposition 4.21 *Let P be a finite point set in a Euclidean space \mathbb{R}^d and $m \in [0, 1[$ be a mass parameter. Considering the empirical measure μ on P ,*

$$d_B^{\log}(\text{Dgm}(d_{\mu,m}), \text{Dgm}(d_{\mu,m}^P)) \leq \ln(\sqrt{3}),$$

$$d_B^{\log}(\text{Dgm}(d_{\mu,m}), \text{Dgm}(d_{\mu,m}^W)) \leq \ln(\sqrt{6}).$$

Proof: Again the proof is a simple use of Theorems 4.7 and 4.15 to obtain an interleaving between sub-level sets. Then, we use Theorem 2.23. ■

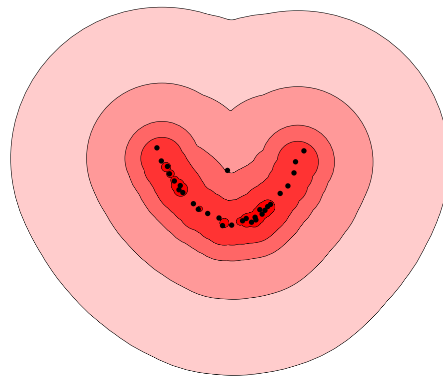
Note that the approximations are not optimal for persistence diagram. A simple rescaling of $d_{\mu,m}^P$ and $d_{\mu,m}^W$ reduces the interleaving factor in logarithmic scale to $\ln(6^{\frac{1}{4}})$. However, we always use the same approximation in practice and the rescaling does not change the shape of the persistence diagram. The only influence is a slight shift of the diagram.

4.6 Example of approximations

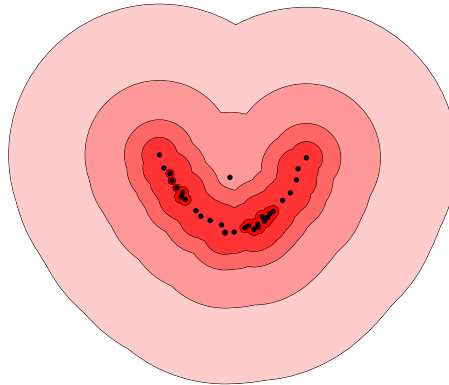
The witnessed k -distance and the approximation by a power distance with support on the points give guarantees of the same quality. However, there are some slight differences in the shape of the sub-level sets. In an Euclidean space \mathbb{R}^d , the distance to a measure can be seen as the lower envelope of a set of paraboloids in the space $\mathbb{R}^d \times \mathbb{R}$. The witnessed k -distance removes some of the paraboloids to reduce the complexity, while the power distance approximation completely changes the location of the focus of the paraboloids. Figure 4.6 shows some sub-level sets for these functions using a noisy point set sampled from a paraboloid and containing an outlier. We are using a mass parameter corresponding to $k = 3$.

The difference in shape does not have a huge influence when using the persistence diagrams. However, one has to be aware of it, especially if one wants to use this kind of function to do reconstruction. From a visual point of view, the witnessed k -distance can be seen as the better approximation because it gives a smoother visualisation. The approximation using the power

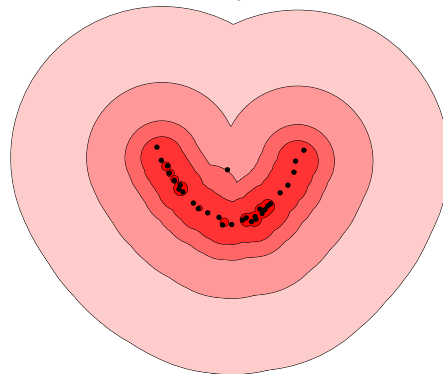
distance supported by P can be used in a more general setting and yields equivalent guarantees for homology inference. The random sampling presents the drawback of being a random algorithm and the need of too many queries to obtain guarantees with high probability, but can be fast in practice and is converging to $d_{\mu,m}$ as the number of queries tends to ∞ .



(a) $d_{\mu,m}$



(b) $d_{\mu,m}^P$



(c) $d_{\mu,m}^W$

Figure 4.6 – Sub-level sets

5 Persistence of power distance functions approximation

In the previous chapter, we showed that distance to a measure functions can be approximated by power distances. We now consider the computation of the persistence diagram for such functions. Methods used to compute persistence diagrams of distance functions, using simplicial complexes such as the Čech complex of the α -complex [56] can be adapted. However, it requires building simplicial complexes of large size that makes it impossible to use for high dimensional data.

In this chapter, we adapt the classic Vietoris-Rips filtration to a weighted setting in order to approximate the persistence diagram of a power distance function. The weighted structure has an interesting structure. It is stable with respect to small perturbations of the support points or the weights, and it almost induces a graph metric on the simplicial complex. This leads to the hope of sparsifying the filtration in order to reduce its complexity, using the method from [88]. Unfortunately, this does not work and we show a way to use the linear-sized approximation without weights to create a linear-sized approximation for the weighted setting.

5.1 Weighted Rips filtration

Given a weighted set (P, w) and the associated power distance f , one can introduce a generalization of the Rips filtration that is adapted to the weighted setting as has been done in [67]. Combined with the approximation of $d_{\mu, m}$ by a power distance, this construction allows us to approximate the persistence diagram of $d_{\mu, m}$. Moreover, we show that it is stable with respect to perturbation of the underlying sample and for similar weighted sets. This gives an interest of its own to the weighted Rips filtration, that can therefore be used as a topological signature.

5.1.1 Definition

The sub-level set $f^{-1}([-\infty, \alpha])$ is the union of the balls centred on the points p of P with radius $r_p(\alpha) = \sqrt{\alpha^2 - w_p^2}$ for $\alpha \geq 0$ and the empty set for $\alpha < 0$. By convention, we consider that the ball is empty when the radius is imaginary and fix $r_p(\alpha) = i$ if $\alpha < 0$. We can define the nerve of this union:

Definition 5.1 *Let (P, w) be a weighted set in a metric space \mathbb{X} , the weighted Čech complex*

$C_\alpha(P, w)$ for parameter α is defined as the union of simplices σ such that $\bigcap_{p \in \sigma} B(p, r_p(\alpha)) \neq \emptyset$.

The Čech complex can be difficult to compute due to the problem of testing if a collection of metric balls has a common intersection. It is often approximated by the Rips complex which requires only distance computations. We do the same here and define a weighted version of the Rips complex.

Definition 5.2 For a weighted set (P, w) in a metric space \mathbb{X} , the weighted Rips complex $R_\alpha(P, w)$ for a parameter α is the maximal simplicial complex whose 1-skeleton has an edge for each pair (p, q) such that $r_p(\alpha)$ and $r_q(\alpha)$ are real and $d_{\mathbb{X}}(p, q) < r_p(\alpha) + r_q(\alpha)$. The weighted Rips filtration is the sequence $\{R_\alpha(P, w)\}$ for all $\alpha \in \mathbb{R}$.

Remark that if all weights are equal to 0, we are in the classical case of balls with equal radii. We use the weighted Rips filtration to approximate the weighted Čech filtration thanks to the following interleaving. For simplicity, the notation (P, w) indicating the point set P with weights w is omitted.

Proposition 5.3 If (P, w) is a weighted set on a metric space \mathbb{X} , then for all $\alpha \in \mathbb{R}$:

$$C_\alpha \subseteq R_\alpha \subseteq C_{2\alpha}.$$

Proof: Let α be a real number. The first inclusion is obtained by the definition of the weighted Rips complex. Let (p, q) be an edge such that $(p, q) \in C_\alpha$. Then $B(p, r_p(\alpha)) \cap B(q, r_q(\alpha)) \neq \emptyset$, which means that $d_{\mathbb{X}}(p, q) < r_p(\alpha) + r_q(\alpha)$, i.e. $(p, q) \in R_\alpha$.

For the other inclusion, let σ be a simplex of R_α . We fix p_0 to be the point of σ with the greatest weight. This means that for any $p \in \sigma$, $r_p(\alpha) \geq r_{p_0}(\alpha)$. Since $\sigma \in R_\alpha$, we get that, for all p and q in σ , $d_{\mathbb{X}}(p, q) < r_p(\alpha) + r_q(\alpha)$ with both radii real. To prove that $\sigma \in C_{2\alpha}$ we need to prove that:

$$\bigcap_{p \in \sigma} B(p, r_p(2\alpha)) \neq \emptyset.$$

We prove that p_0 belongs to this intersection. For each $p \in \sigma$:

$$d_{\mathbb{X}}(p, p_0) < r_p(\alpha) + r_{p_0}(\alpha) \leq 2 r_p(\alpha) = \sqrt{(2\alpha)^2 - 4w_p^2} \leq r_p(2\alpha).$$

Hence $p_0 \in B(p, r_p(2\alpha))$. ■

5.1.2 Stability

The persistence diagram of a weighted Rips filtration $\{R_\alpha(P, w)\}$ is stable under small perturbations of the set P . It can thus be used in applications like signatures in the spirit of [23]. In order to speak about the persistence diagram of a weighted Rips filtration, we first need to verify that the filtration is q -tame. This is always the case when the set P is compact. Remark that we are not restricted to finite point sets. We first consider the persistence modules associated with the weighted Rips filtrations and show that these modules are closely interleaved when the sets P and Q are close.

Lemma 5.4 *Let P, Q be two subsets of a metric space \mathbb{X} and let $w : \mathbb{X} \rightarrow \mathbb{R}$ be a t -Lipschitz function. Then $H_*(\{R_\alpha(P, w)\})$ and $H_*(\{R_\alpha(Q, w)\})$ are ϵ -interleaved for $\epsilon = \sqrt{1+t^2}d_H(P, Q)$.*

Proof: We need to show that there exists ϵ -homomorphisms π_{P*} and π_{Q*} such that $\pi_{P*}\pi_{Q*} = 1_{H_*(R_\alpha(P, w))}^{2\epsilon}$ and $\pi_{Q*}\pi_{P*} = 1_{H_*(R_\alpha(Q, w))}^{2\epsilon}$. To do so, we need three steps. First, we build simplicial maps $R_\alpha(P, w) \rightarrow R_{\alpha+\epsilon}(Q, w)$ and $R_\alpha(Q, w) \rightarrow R_{\alpha+\epsilon}(P, w)$ for every α . The construction is not necessarily unique but the simplicial maps are contiguous and hence induce the same homomorphism. Finally, we show that the simplicial maps in Figure 5.1 are contiguous and thus the persistence modules are ϵ -interleaved.

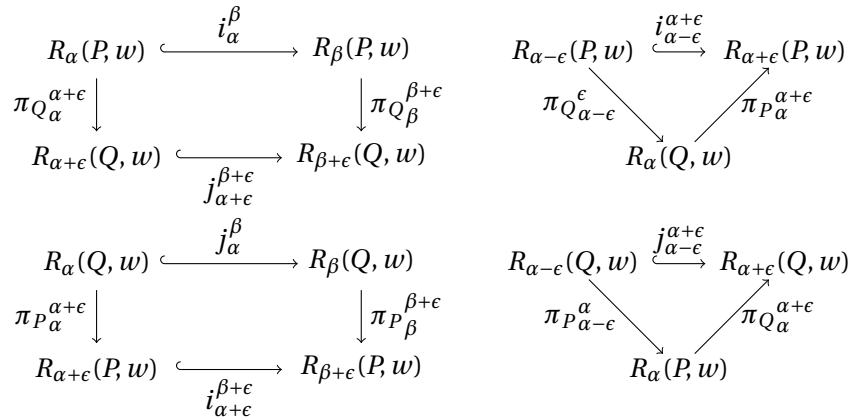


Figure 5.1 – Diagrams with contiguous simplicial maps between Rips filtrations

The simplicial maps $i_\alpha^\beta : R_\alpha(P, w) \rightarrow R_\beta(P, w)$ and $j_\alpha^\beta : R_\alpha(Q, w) \rightarrow R_\beta(Q, w)$ for $\alpha < \beta$ are the canonical inclusions. We consider two maps $\pi_P : Q \rightarrow P$ and $\pi_Q : P \rightarrow Q$ such that $d_{\mathbb{X}}(p, \pi_Q(p)) \leq d_H(P, Q)$ and $d_{\mathbb{X}}(q, \pi_P(q)) \leq d_H(P, Q)$ for any $p \in P$ and $q \in Q$. By definition of the Hausdorff distances, such maps always exist¹. Let us show that these maps induce simplicial maps.

Let us consider the function π_P and let us fix $\alpha > 0$. Let (q', q'') be an edge of $R_\alpha(Q, w)$. It means that $B(q', r_{q'}(\alpha)) \cap B(q'', r_{q''}(\alpha)) \neq \emptyset$. Lemma 4.9 implies that for any $q \in Q$, $B(q, r_q(\alpha)) \subset B(\pi_P(q), r_{\pi_P(q)}(\alpha + \sqrt{1+t^2}d_H(P, Q)))$. Thus, $(\pi_P(q'), \pi_P(q''))$ is an edge of $R_{\alpha+\epsilon}(P, w)$ because:

$$B(\pi_P(q'), r_{\pi_P(q')}(\alpha + \epsilon)) \cap B(\pi_P(q''), r_{\pi_P(q'')}(\alpha + \epsilon)) \supset B(q', r_{q'}(\alpha)) \cap B(q'', r_{q''}(\alpha)) \neq \emptyset.$$

As $R_\alpha(P, w)$ is a clique complex for any α , this is sufficient to prove that π_P induce a family of simplicial maps $\{\pi_{P_\alpha}^{\alpha+\epsilon}\}$. The roles of P and Q are symmetric. Therefore, the result holds for π_Q as well.

Notice that the construction of π_P and π_Q is not unique. However, we now show that if π_P and π'_P are simplicial maps verifying $d_{\mathbb{X}}(q, \pi_P(q)) \leq d_H(P, Q)$ and $d_{\mathbb{X}}(q, \pi'_P(q)) \leq d_H(P, Q)$ for

¹Strictly speaking, such points may not exist if P or Q are not compact. However, we can replace $d_H(P, Q)$ by any $\eta > d_H(P, Q)$ in the whole proof and finally go to the limit to obtain the same result.

any $q \in Q$, then they are contiguous. π_P therefore induce a canonical homomorphism at the homology level thanks to Theorem 2.11.

Let α be a parameter and (q, q') be an edge of $R_\alpha(Q, w)$. By definition, $B(q, r_q(\alpha)) \cap B(q', r_{q'}(\alpha)) \neq \emptyset$. Moreover, using Lemma 4.9,

$$B(\pi_P(q), r_{\pi_P(q)}(\alpha + \epsilon)) \cap B(\pi'_P(q), r_{\pi'_P(q)}(\alpha + \epsilon)) \supset B(q, r_q(\alpha))$$

and thus

$$\begin{aligned} & B(\pi_P(q), r_{\pi_P(q)}(\alpha + \epsilon)) \cap B(\pi'_P(q), r_{\pi'_P(q)}(\alpha + \epsilon)) \\ & \cap B(\pi_P(q'), r_{\pi_P(q')}(\alpha + \epsilon)) \cap B(\pi'_P(q'), r_{\pi'_P(q')}(\alpha + \epsilon)) \neq \emptyset. \end{aligned}$$

Hence the simplex generated by $\{\pi_P(q), \pi_P(q'), \pi'_P(q), \pi'_P(q')\}$ belongs to the complex $C_{\alpha+\epsilon}(P, w)$. A fortiori it belongs to $R_{\alpha+\epsilon}(P, w)$ and Lemma 2.12 implies the contiguity of π_P and π'_P . Again, the same can be applied after exchanging P and Q .

Thus we can choose two arbitrary projections π_P and π_Q as they all induce the same canonical homomorphisms π_{P*} and π_{Q*} . To prove that $\pi_{P*}\pi_{Q*} = 1_{H_*(R_\alpha(P, w))}^{2\epsilon}$, we now prove that the diagrams from Figure 5.1 commute with contiguous maps. Taking advantage of the symmetry of the problem, we only prove the two diagrams from the first line.

Let us fix $\alpha < \beta$. The first diagram commutes if and only if $\pi_{Q\beta}^{\beta+\epsilon} \circ Id_{P\alpha}^\beta$ and $Id_{Q\alpha+\epsilon}^{\beta+\epsilon} \circ \pi_{Q\alpha}^{\alpha+\epsilon}$ are contiguous. The two functions are equal and thus contiguous because $\pi_{Q\alpha}^{\alpha+\epsilon}$ and $\pi_{Q\beta}^{\beta+\epsilon}$ are induced by the same function π_Q .

Let us now prove that $\pi_{P\alpha}^{\alpha+\epsilon} \circ \pi_{Q\alpha-\epsilon}^\alpha$ and $i_{\alpha-\epsilon}^{\alpha+\epsilon}$ are contiguous for any α . Let us fix α and let (p, p') be an edge of $R_{\alpha-\epsilon}(P, w)$. By definition, $B(p, r_p(\alpha - \epsilon)) \cap B(p', r_{p'}(\alpha - \epsilon)) \neq \emptyset$. Moreover, using Lemma 4.9 we get:

$$B(p, r_p(\alpha - \epsilon)) \subset B(\pi_Q(p), r_{\pi_Q(p)}(\alpha)) \subset B(\pi_P \circ \pi_Q(p), r_{\pi_P \circ \pi_Q(p)}(\alpha + \epsilon)).$$

The same holds for p' and thus:

$$\begin{aligned} & B(p, r_p(\alpha + \epsilon)) \cap B(\pi_P \circ \pi_Q(p), r_{\pi_P \circ \pi_Q(p)}(\alpha + \epsilon)) \\ & \cap B(p', r_{p'}(\alpha + \epsilon)) \cap B(\pi_P \circ \pi_Q(p'), r_{\pi_P \circ \pi_Q(p')}(\alpha + \epsilon)) \neq \emptyset. \end{aligned}$$

Thus the simplex generated by $\{i_{\alpha-\epsilon}^{\alpha+\epsilon}(p), i_{\alpha-\epsilon}^{\alpha+\epsilon}(p'), \pi_{P\alpha}^{\alpha+\epsilon} \circ \pi_{Q\alpha-\epsilon}^\alpha(p), \pi_{P\alpha}^{\alpha+\epsilon} \circ \pi_{Q\alpha-\epsilon}^\alpha(p')\}$ is in $C_{\alpha+\epsilon}(P, w) \subset R_{\alpha+\epsilon}(P, w)$. Lemma 2.12 guarantees that $\pi_{P\alpha}^{\alpha+\epsilon} \circ \pi_{Q\alpha-\epsilon}^\alpha$ and $i_{\alpha-\epsilon}^{\alpha+\epsilon}$ are contiguous. From before, $\{\pi_{P\alpha}^{\alpha+\epsilon} \circ \pi_{Q\alpha-\epsilon}^\alpha\}$ induces the 2ϵ -homomorphism $\pi_{P*}\pi_{Q*}$. By definition, $\{i_{\alpha-\epsilon}^{\alpha+\epsilon}\}$ induces $1_{H_*(R_\alpha(P, w))}^{2\epsilon}$. Using Theorem 2.11, we have $\pi_{P*}\pi_{Q*} = 1_{H_*(R_\alpha(P, w))}^{2\epsilon}$.

By symmetry of the roles of P and Q , $\{R_\alpha(P, w)\}$ and $\{R_\alpha(Q, w)\}$ are ϵ -interleaved. ■

From this lemma, we derive that the weighted Rips filtration associated with a pair (P, w) , where P is compact, is q -tame and hence has a well defined persistence diagram.

Proposition 5.5 *Let P be a subset of a metric space \mathbb{X} and let $w : \mathbb{X} \rightarrow \mathbb{R}$ be a t -Lipschitz function. If P is compact, then $\{R_\alpha(P, w)\}_{\alpha \in \mathbb{R}}$ is q -tame.*

Proof: We show that, for any $\epsilon > 0$, we can build a finite persistence module which is ϵ -interleaved with the persistence module of $\{R_\alpha(P, w)\}$. A finite persistence module is a fortiori locally finite and [26, Theorem 4.19] induces the q -tameness of $\{R_\alpha(P, w)\}$.

Fix $\epsilon > 0$. P is compact. As a consequence, there exists a finite point set Q of P such that $d_H(P, Q) \leq \frac{\epsilon}{\sqrt{1+t^2}}$. The persistence module of $\{R_\alpha(Q, w)\}$ is finite and therefore locally finite. Moreover, using Lemma 5.4, $\{R_\alpha(Q, w)\}$ and $\{R_\alpha(P, w)\}$ are ϵ -interleaved. Hence $\{R_\alpha(P, w)\}$ is q -tame using [26, Theorem 4.19]. ■

Talking about different filtrations, we can relate the persistence diagrams of the weighted Rips filtration through the following stability theorem.

Theorem 5.6 *Let P and Q be two compact subsets of a metric space \mathbb{X} . Let $w : \mathbb{X} \rightarrow \mathbb{R}$ be a t -Lipschitz function. Then,*

$$d_B(\text{Dgm}(\{R_\alpha(P, w)\}), \text{Dgm}(\{R_\alpha(Q, w)\})) \leq \sqrt{1+t^2} d_H(P, Q).$$

Proof: P and Q are two compact sets and thus the diagrams are well-defined thanks to Proposition 5.5 that guarantees the q -tameness of the filtrations. Lemma 5.4 implies that $H_*(\{R_\alpha(P, w)\})$ and $H_*(\{R_\alpha(Q, w)\})$ are $\sqrt{1+t^2} d_H(P, Q)$ -interleaved. The relation between the persistence diagrams is obtained by applying Theorem 2.23. ■

The embedding of both P and Q in the same space is not required to work with the Rips filtrations. It is possible to compare the diagrams using the notion of ϵ -correspondence from [27].

Definition 5.7 *Let P and Q be two points sets in metric spaces \mathbb{X} and \mathbb{X}' . Let C be a subset of $P \times Q$. C is an ϵ -correspondence if:*

$$\forall p \in P, \exists q \in Q, (p, q) \in C$$

$$\forall q \in Q, \exists p \in P, (p, q) \in C$$

$$\forall p, p' \in P, \forall q, q' \in Q, (p, q) \in C \wedge (p', q') \in C \implies |d_{\mathbb{X}}(p, p') - d_{\mathbb{X}'}(q, q')| \leq \epsilon$$

Theorem 5.8 *Let \mathbb{X} and \mathbb{X}' be two metric spaces, P and Q be two finite point sets of \mathbb{X} and \mathbb{X}' respectively. Let C be an ϵ -correspondence between P and Q . Let $w_{\mathbb{X}} : \mathbb{X} \rightarrow \mathbb{R}$ and $w_{\mathbb{X}'} : \mathbb{X}' \rightarrow \mathbb{R}$ be two t -Lipschitz functions such that if $(p, q) \in C$ then $|w_{\mathbb{X}}(p) - w_{\mathbb{X}'}(q)| \leq t\epsilon$. Then:*

$$d_B(\text{Dgm}(\{R_\alpha(P, w_{\mathbb{X}})\}), \text{Dgm}(\{R_\alpha(Q, w_{\mathbb{X}'})\})) \leq (1+t)\epsilon$$

Proof: Let us remark that the proof of Lemma 5.4 is valid when considering correspondences. We can build maps π_Q and π_P in the same way. For any $p \in P$, we choose π_Q among the q such that $(p, q) \in C$. All the different maps we can build this way are contiguous, adapting the proof of Lemma 5.4 and thus the construction is canonical at the homology level. Moreover, the proof of interleaving can be adapted to get that $H_*(\{R_\alpha(P, w_{\mathbb{X}})\})$ and $H_*(\{R_\alpha(Q, w_{\mathbb{X}'})\})$ are ϵ -interleaved. Proposition 5.5 guarantees the tameness of the filtrations and the relation between persistence diagrams is given by Theorem 2.23. ■

Persistence diagrams can thus be used as signatures on a wider class of objects. This makes the weighted Rips filtration useful for other application than the computation of approximate persistence diagrams.

5.1.3 Approximation of $d_{\mu,m}$

To use the weighted Rips filtration to approximate the persistence diagram of the distance to a measure, we need to restrict the class of spaces considered. If the intersection of any finite number of balls in \mathbb{X} is either contractible or empty, \mathbb{X} is said to have the *good cover property*. Then the Nerve Theorem [69, 4G.3] guarantees that the Čech complex, which is the nerve of a union of balls, has the same homology as this union. We can also compute the persistence diagram thanks to the Persistent Nerve Lemma [35]. We obtain an approximation of $\text{Dgm}(d_{\mu,m})$ using the weighted Rips filtration. Remark that not all metric spaces, like the sphere S^d , have the good cover property.

Theorem 5.9 *Let \mathbb{X} be a triangulable metric space with the good cover property and let P be a finite point set of \mathbb{X} . Considering the empirical measure μ over P and $m \in [0, 1[$ a mass parameter, we obtain on a logarithmic scale:*

$$d_B^{\log}(\text{Dgm}(d_{\mu,m}), \text{Dgm}(\{R_\alpha(P, d_{\mu,m}^P)\})) \leq \ln(2\sqrt{5}).$$

Proof: Given that \mathbb{X} is triangulable, the sub-level sets filtration of $d_{\mu,m}$ is q -tame by Proposition 4.18. The persistence diagram $\text{Dgm}(d_{\mu,m})$ is thus well-defined. Recall that $d_{\mu,m}$ is a 1-Lipschitz function from Corollary 3.16. P is a compact subset of \mathbb{X} and therefore $\text{Dgm}(R_\alpha(P, d_{\mu,m}^P))$ is well-defined according to Proposition 5.5.

We approximate $d_{\mu,m}$ with $d_{\mu,m}^P$. The result of Theorem 4.12 gives a $\sqrt{5}$ multiplicative inter-leaving. For any $\alpha \in \mathbb{R}$,

$$d_{\mu_P,m}([-\infty, \alpha]) \subset d_{\mu_P,m}^P([-\infty, \sqrt{2}\alpha]) \subset d_{\mu_P,m}([-\infty, \sqrt{10} d_{\mu_P,m}^P]).$$

So, Theorem 2.23 implies

$$d_B^{\log}(\text{Dgm}(d_{\mu,m}), \text{Dgm}(d_{\mu,m}^P)) \leq \ln(\sqrt{5}).$$

By the Persistent Nerve Lemma, the sub-level sets filtration of $d_{\mu,m}^P$, which is a union of balls, has the same persistent homology as the filtration of its nerve, $\{C_\alpha(P, w)\}$. Thus, Proposition 5.3 and Theorem 2.23 imply

$$d_B^{\log}(\text{Dgm}(d_{\mu,m}^P), \text{Dgm}(\{R_\alpha(P, d_{\mu,m}^P)\})) \leq \ln(2).$$

The triangle inequality for the bottleneck distance gives the desired result. ■

5.2 Weighted Rips induced metric

Considering a weighted Rips filtration, we look at the apparition time of vertices and edges. As a Rips filtration is a clique complex, this completely defines the filtration. When the parameter

of the filtration tend to $+\infty$, all possible edges belongs to the filtration. We can use the parameter at which they first appear as a pseudo-distance.

Definition 5.10 Let (P, w) be a weighted point set in a metric space \mathbb{X} . Given $(p, q) \in P^2$, we define the pseudo-distance f :

$$f(p, q) = \begin{cases} \sqrt{\frac{w_p^2 + w_q^2}{2} + \frac{d_{\mathbb{X}}(p, q)^2}{4} + \frac{(w_p^2 - w_q^2)^2}{4d_{\mathbb{X}}(p, q)^2}} & \text{if } |w_p^2 - w_q^2| < d_{\mathbb{X}}(p, q)^2 \\ \max(w_p, w_q) & \text{otherwise} \end{cases}$$

Proposition 5.11 Let (P, w) be a weighted point set in a metric space \mathbb{X} . Given $(p, q) \in P^2$ and $\alpha \geq 0$, the edge (p, q) belongs to $R_\alpha(P, w)$ if and only if $f(p, q) < \alpha$.

Proof: We are looking for the first time α the two weighted balls around p and q intersect. Recall that their radii for a parameter α is given by $r_p(\alpha) = \sqrt{\alpha^2 - w_p^2}$ if $\alpha \geq w_p$ and the ball is empty otherwise. If $p = q$, then the balls intersect as soon as they exists, i.e. for $\alpha = w_p$. Let us consider the case $p \neq q$. The radii are continuous with respect to the parameter α . Thus, there exist two ways for balls to intersect for the first time as shown in Figure 5.2. Either the two balls grow and become tangent or one of them is empty until after its centre is covered by the other one.

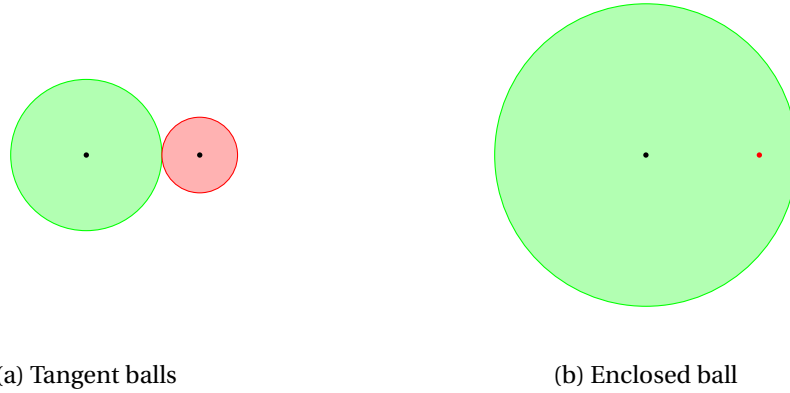


Figure 5.2 – The two kind of first intersection between balls

Let us start with the second case. Assuming without loss of generality that $w_p < w_q$, it means that the ball centred in q appears at a time when the ball centred in p contains q . Id est, there exists an α such that $w_q \geq \alpha$ and $r_p(\alpha) > d_{\mathbb{X}}(p, q)$. In other words $w_q^2 - w_p^2 > d_{\mathbb{X}}(p, q)^2$. In this case the two balls intersect as soon as the ball centred in q appears, which happens for $\alpha = w_q$. Relaxing the assumption $w_q > w_p$, this implies that if $|w_q^2 - w_p^2| > d_{\mathbb{X}}(p, q)^2$, then the balls intersect for $\alpha \geq \max(w_p, w_q)$.

Let us get back to the case of tangent balls. Two balls will intersect for the first time by being tangent for parameter α such that $\alpha \geq w_p$ and $\alpha \geq w_q$. The relation between radii and distance

can be written as $r_p(\alpha) + r_q(\alpha) = d_{\mathbb{X}}(p, q)$.

$$\begin{aligned} d_{\mathbb{X}}(p, q) &= \sqrt{\alpha^2 - w_p^2} + \sqrt{\alpha^2 - w_q^2} \\ d_{\mathbb{X}}(p, q)^2 &= 2\alpha^2 - (w_p^2 + w_q^2) + 2\sqrt{\alpha^4 - \alpha^2(w_p^2 + w_q^2) + w_p^2 w_q^2} \\ (d_{\mathbb{X}}(p, q)^2 - 2\alpha^2 + (w_p^2 + w_q^2))^2 &= 4(\alpha^4 - \alpha^2(w_p^2 + w_q^2) + w_p^2 w_q^2) \end{aligned}$$

$$\left(d_{\mathbb{X}}(p, q)^2 + (w_p^2 + w_q^2)\right)^2 + 4\alpha^4 - 4\alpha^2(d_{\mathbb{X}}(p, q)^2 + w_p^2 + w_q^2) = 4\alpha^4 - 4\alpha^2(w_p^2 + w_q^2) + 4w_p^2 w_q^2$$

$$\begin{aligned} 4\alpha^2 d_{\mathbb{X}}(p, q)^2 &= d_{\mathbb{X}}(p, q)^4 + (w_p^2 + w_q^2)^2 + 2d_{\mathbb{X}}(p, q)^2(w_p^2 + w_q^2) - 4w_p^2 w_q^2 \\ 4\alpha^2 d_{\mathbb{X}}(p, q)^2 &= d_{\mathbb{X}}(p, q)^4 + 2d_{\mathbb{X}}(p, q)^2(w_p^2 + w_q^2) + (w_p^2 - w_q^2)^2 \\ \alpha^2 &= \frac{d_{\mathbb{X}}(p, q)^2}{4} + \frac{(w_p^2 + w_q^2)^2}{2} + \frac{(w_p^2 - w_q^2)^2}{4d_{\mathbb{X}}(p, q)^2} \end{aligned}$$

Remark that, assuming $w_p \geq w_q$, if $d_{\mathbb{X}}(p, q)^2 = |w_p^2 - w_q^2| \neq 0$, then

$$\sqrt{\frac{w_p^2}{2} + \frac{w_q^2}{2} + \frac{d_{\mathbb{X}}(p, q)^2}{4} + \frac{(w_p^2 - w_q^2)^2}{4d_{\mathbb{X}}(p, q)^2}} = w_p.$$

Hence Definition 5.10 is coherent.

Finally, consider $(p, q) \in P^2$. Then $(p, q)^2 \in \mathbb{R}_{\alpha}(P, w)$ if and only if $B(p, r_p(\alpha)) \cap B(q, r_q(\alpha)) \neq \emptyset$, which is equivalent to $f(p, q) < \alpha$. ■

The expression of the time of first intersection gives an easy way to compute the apparition time of edges in the weighted Rips filtration. In addition to its interest for computing the filtration, it provides the function f that has some interesting properties. It is almost a metric on the full Rips complex. The function is symmetric and positive but does not possess the relation $f(p, p) = 0$. In fact, $p \neq q$ implies that $f(p, q) > 0$ but the other direction does not hold. Two approaches are then possible to use it. Either, we can force $f(p, p)$ to be equal to 0 or we can use directly f to interpolate a metric on the complex. In both cases, we need f to respects the triangle inequality.

Theorem 5.12 *Let (P, w) be a weighted point set and f the function induced by the weighted Rips filtration. Then,*

$$\forall a, b, c \in P, f(a, c) \leq f(a, b) + f(b, c)$$

The proof, rather long and purely technical, is detailed in Appendix A.

5.3 Sparse weighted Rips

The weighted Rips filtration has the desired approximation guarantees to be used for persistence computation, but like the classic Rips filtration for points, it usually becomes too large to be computed in full. Several approaches have been proposed to reduce this complexity in

the case without weights. One can try to simplify the simplicial complexes while guaranteeing its topology [6] or looks for a complex that approximates closely the Rips complex [51, 88]. Here, we consider a construction of the latter form.

In [88], it was shown how to construct a filtration $\{S_\alpha\}$, called the *sparse Rips filtration*, that gives a provably good approximation to the Rips filtration and has size linear in the number of points for metrics with constant doubling dimension. Specifically, for a user-defined parameter ϵ , the log-bottleneck distance between the persistence diagrams of the Sparse Rips filtration and the Rips filtration is at most ϵ . The goal of this section is to extend that result to weighted Rips filtrations.

The sparsification technique cannot be used directly here, since the power distance does not exactly induce a metric. The fact that points does not appear immediately in the filtration, expressed by the fact that the induced function $f(p, p)$ can be non zero, can create phenomenon where the weighted point set has a greater doubling dimension than the point set without weights. For example, consider the case of points regularly located on a cycle. The intrinsic dimension of the object is 1 and thus the sparse Rips filtration has a size linear in the number of points n times C where C is a constant depending on ϵ . If the intrinsic dimension, the size is $O(C^l n)$. However, if we set weights at all points to be some large constant, then all points are now at the same distance from each other and the doubling dimension becomes $\log n$. Thus the size of the sparse Rips filtration applied on the weighted setting will be quadratic in n .

In this chapter, we show that it is possible to approximate the weighted Rips filtration by using the sparse Rips filtration for the case without weights and changing the time of apparition of simplices. This guarantees that the size will be the same as the one from [88] while approximating the weighted Rips filtration. For the rest of this section, we fix a weighted point set (P, w) in a metric space \mathbb{X} , where the weight function $w : \mathbb{X} \rightarrow \mathbb{R}$ is t -Lipschitz, for some constant t . To simplify notation, we let R_α denote the weighted Rips complex $R_\alpha(P, w)$.

Definition 5.13 *Given a weighted point set (P, w) in a metric space \mathbb{X} , the sparse weighted Rips filtration $\{T_\alpha\}$ of (P, w) is defined as*

$$\forall, \alpha \in \mathbb{R}, T_\alpha = S_\alpha \cap R_\alpha,$$

where $\{S_\alpha\}$ is the sparse Rips filtration of P and $\{R_\alpha\}$ is the weighted Rips filtration of (P, w) .

The sparse Rips filtration $\{S_\alpha\}$ captures the underlying metric space and the weighted Rips filtration $\{R_\alpha\}$ captures the structure of the sub-level sets of the power distance function. Computing $\{T_\alpha\}$ can be done in $O(n^2)$ by first computing $\{S_\alpha\}$ and then reordering the simplices according to the birth time in $\{R_\alpha\}$. This is equivalent to filtering the complex S_∞ . Note that the sparsification depends only on the metric, and not on the weights. Thus, the same sparse Rips complex can be used as the underlying complex for different weight functions.

The technical challenge is to relate the persistence diagram of this new filtration to the persistence diagram of the weighted Rips filtration as in the following theorem.

Theorem 5.14 *Let (P, w) , be a finite, weighted subset of a metric space \mathbb{X} with t -Lipschitz weights. Let $\epsilon \in]0, 1[$ be a fixed constant used in the construction of the sparse weighted Rips*

filtration $\{T_\alpha\}$. Then,

$$d_B^{\log}(\text{Dgm}(\{T_\alpha\}), \text{Dgm}(\{R_\alpha\})) \leq \ln \left(\frac{1 + \sqrt{1 + t^2} \epsilon}{1 - \epsilon} \right).$$

Since these filtrations are not interleaved, the only hope is to find an interleaving of the persistence modules, which requires finding suitable homomorphisms between the homology groups of the different filtrations. After detailing a new construction of the sparse Rips filtration from [88], which uses a furthest point sampling instead of the original net tree structure, the rest of this section proves Theorem 5.14.

5.3.1 Sparse Rips complexes

Let (p_1, \dots, p_n) be a greedy permutation of the points P in a finite metric space \mathbb{X} . That is, p_1 is chosen arbitrarily and $p_i = \arg\max_{p \in P \setminus P_{i-1}} d_{\mathbb{X}}(p, P_{i-1})$, where $P_{i-1} = \{p_1, \dots, p_{i-1}\}$ is the $(i-1)^{\text{st}}$ prefix. We define the *insertion radius* λ_{p_i} of point p_i to be

$$\lambda_{p_i} = d_{\mathbb{X}}(p_i, P_{i-1}).$$

To avoid excessive subscripts, we write λ_i in place of λ_{p_i} when we know the index of p_i . We adopt the convention that $\lambda_1 = \infty$ and $\lambda_{n+1} = 0$. The greedy permutation has the nice property that each prefix P_i is a λ_i -net in the sense that

1. $d_{\mathbb{X}}(p, P_i) \leq \lambda_i$ for all $p \in P$.
2. $d_{\mathbb{X}}(p, q) \geq \lambda_i$ for all $p, q \in P_i$.

We extend these nets to an arbitrary parameter γ :

$$\begin{aligned} N_\gamma &= \{p \in P \mid \lambda_p > \gamma\}. \\ \overline{N}_\gamma &= \{p \in P \mid \lambda_p \geq \gamma\}. \end{aligned}$$

Note that for all $p \in P$, $d_{\mathbb{X}}(p, N_\gamma) \leq \gamma$ and $d_{\mathbb{X}}(p, \overline{N}_\gamma) < \gamma$.

One way to get a sparse Rips-like filtration is to take a union of Rips complexes on the nets N_γ . However, this can add noise to the persistence diagram compared to the Rips filtrations as illustrated in Figure 5.3. It shows two different sub-level sets of the distance to the points. On the left, there is no sparsification. As soon as the four balls form a connected set, the homology is trivial for the rest of their growth. On the right, after the four balls connect, we sparsify naively by removing the central point, which is covered by the other balls. However, the light sub-level set has now a non-trivial homology in dimension 1.

This phenomenon can be avoided by a careful perturbation of the distance. For a point p , the perturbation varies with the scale and is defined as follows:

$$s_p(\alpha) = \begin{cases} 0 & \text{if } \alpha \leq \frac{\lambda_p}{\epsilon} \\ \alpha - \frac{\lambda_p}{\epsilon} & \text{if } \frac{\lambda_p}{\epsilon} < \alpha < \frac{\lambda_p}{\epsilon(1-\epsilon)} \\ \epsilon \alpha & \text{if } \frac{\lambda_p}{\epsilon(1-\epsilon)} \leq \alpha \end{cases}$$

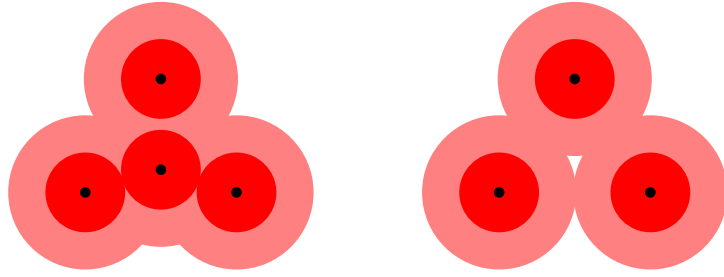
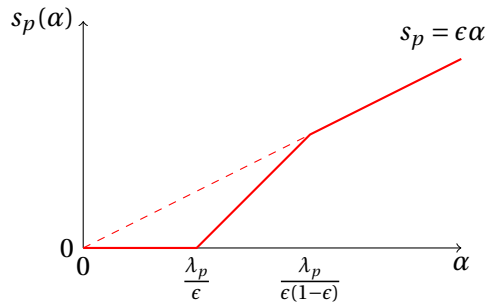


Figure 5.3 – Topological noise induced by naive sparsification


 Figure 5.4 – Perturbation in function of α

Note that s_p is 1-Lipschitz. The resulting perturbed distance is defined as

$$f_\alpha(p, q) = d_\mathbb{X}(p, q) + s_p(\alpha) + s_q(\alpha).$$

For any fixed p and q , the Lipschitz property of s_p and s_q implies that for all $\alpha \leq \beta$:

$$f_\alpha(p, q) \leq f_\beta(p, q) + 2(\beta - \alpha).$$

Definition 5.15 Given the nets N_γ and the distance function f_α , we define the sparse Rips complex at scale α as

$$Q_\alpha = \{\sigma \subset \overline{N}_{\epsilon(1-\epsilon)\alpha} \mid \forall p, q \in \sigma, f_\alpha(p, q) < 2\alpha\}.$$

On its own, the sequence of complexes $\{Q_\alpha\}$ does not form a filtration. Some simplices disappear as α grows.

Definition 5.16 The sparse Rips filtration is defined as:

$$S_\beta = \bigcup_{\alpha \leq \beta} Q_\alpha.$$

5.3.2 Projection onto Nets

To relate sparse Rips complexes with Rips complexes, we build a collection of projections of the points onto the nets.

$$\pi_\alpha(p) = \begin{cases} p & \text{if } p \in N_{\epsilon(1-\epsilon)\alpha} \\ \operatorname{argmin}_{q \in N_{\epsilon\alpha}} d_{\mathbb{X}}(p, q) & \text{otherwise} \end{cases}$$

For any scale α , the projection π_α maps the points of P to the net $N_{\epsilon(1-\epsilon)\alpha}$. Note that π_α is specifically defined to be a retraction onto $N_{\epsilon(1-\epsilon)\alpha}$. One could wish to use f_α as the distance to define the projection, but it is not known if this yields a retraction.

We present three lemmas about the perturbed distance functions and projections. The projections are used extensively to induce maps between simplicial complexes.

First, we prove that edges do not disappear as the filtration grows.

Lemma 5.17 *If $f_\alpha(p, q) < 2\alpha \leq 2\beta$ then $f_\beta(p, q) < 2\beta$.*

Proof: The proof follows from the definitions f_α and f_β , the Lipschitz property of the perturbations s_p and s_q , and the hypothesis.

$$\begin{aligned} f_\beta(p, q) &= d_{\mathbb{X}}(p, q) + s_p(\beta) + s_q(\beta) \\ &\leq d_{\mathbb{X}}(p, q) + s_p(\alpha) + (\beta - \alpha) + s_q(\alpha) + (\beta - \alpha) \\ &= f_\alpha(p, q) + 2(\beta - \alpha) \\ &< 2\alpha + 2(\beta - \alpha) \\ &= 2\beta. \end{aligned}$$

■

Next, we show that the distance between a point and its projection is at most the change in the perturbed distance.

Lemma 5.18 *For all $q \in P$, $d_{\mathbb{X}}(q, \pi_\alpha(q)) \leq s_q(\alpha) - s_{\pi_\alpha(q)}(\alpha)$, and in particular, $d_{\mathbb{X}}(q, \pi_\alpha(q)) \leq \epsilon\alpha$.*

Proof: Both statements are trivial if $q \in N_{\epsilon(1-\epsilon)\alpha}$, because that would imply that $\pi_\alpha(q) = q$. So, we may assume that $\pi_\alpha(q)$ is the nearest point to q in $N_{\epsilon\alpha}$. It follows that

$$d_{\mathbb{X}}(q, \pi_\alpha(q)) \leq \epsilon\alpha.$$

Moreover, $\lambda_q \leq \epsilon(1-\epsilon)\alpha$, and thus $s_q(\alpha) = \epsilon\alpha$. Also, since $\pi_\alpha(q) \in N_{\epsilon\alpha}$, it must be that $\lambda_{\pi_\alpha(q)} > \epsilon\alpha$ and so $s_{\pi_\alpha(q)} = 0$. Combining these statements, we get

$$d_{\mathbb{X}}(q, \pi_\alpha(q)) \leq \epsilon\alpha = s_q(\alpha) - s_{\pi_\alpha(q)}(\alpha).$$

■

Now, we prove that replacing a point with its projection does not increase the perturbed distance.

Lemma 5.19 For all $p, q \in P$ and all $\alpha \geq 0$, $f_\alpha(p, \pi_\alpha(q)) \leq f_\alpha(p, q)$.

Proof: The statement follows from the definition of f_α , the triangle inequality and Lemma 5.18.

$$\begin{aligned} f_\alpha(p, \pi_\alpha(q)) &= d_{\mathbb{X}}(p, \pi_\alpha(q)) + s_p(\alpha) + s_{\pi_\alpha(q)}(\alpha) \\ &\leq d_{\mathbb{X}}(p, q) + d_{\mathbb{X}}(q, \pi_\alpha(q)) + s_p(\alpha) + s_{\pi_\alpha(q)}(\alpha) \\ &\leq d_{\mathbb{X}}(p, q) + s_p(\alpha) + s_q(\alpha) \\ &= f_\alpha(p, q). \end{aligned}$$

■

5.3.3 Sometimes the projections induce contiguous simplicial maps

In this section, we consider the maps between simplicial complexes that are induced by the projection functions π_α . We are most interested in the case when a pair of projections π_α and π_β induces contiguous simplicial maps between sparse Rips complexes (Lemma 5.22) or weighted Rips complexes (Lemma 5.23). First, we need a couple lemmas that describe the effect of different projections on the endpoints of an edge in sparse or weighted Rips complexes.

Lemma 5.20 Let α, β, γ , and i be such that $\frac{\lambda_{i+1}}{\epsilon(1-\epsilon)} \leq \alpha \leq \beta \leq \gamma \leq \frac{\lambda_i}{\epsilon(1-\epsilon)}$. If an edge (p, q) is in Q_ρ for some $\rho \leq \gamma$ then the edge $(\pi_\alpha(p), \pi_\beta(q)) \in Q_\gamma$.

Proof: The conditions on α, β, γ , and i imply that $\pi_\alpha(p)$ and $\pi_\beta(q)$ are in $\overline{N}_{\epsilon(1-\epsilon)\gamma}$, which is the vertex set of Q_γ .

$$\begin{aligned} \pi_\alpha(p) \in N_{\epsilon(1-\epsilon)\alpha} &\implies \lambda_{\pi_\alpha(p)} > \epsilon(1-\epsilon)\alpha \geq \lambda_{i+1} \\ &\implies \lambda_{\pi_\alpha(p)} \geq \lambda_i \geq \epsilon(1-\epsilon)\gamma \\ &\implies \pi_\alpha(p) \in N_{\epsilon(1-\epsilon)\gamma} \end{aligned}$$

The same holds for $\pi_\beta(q)$, and hence it will suffice to prove that $f_\gamma(\pi_\alpha(p), \pi_\beta(q)) < 2\gamma$ given that $f_\rho(p, q) < 2\rho$. Next we consider three cases depending on the value of ρ in relation to α and β .

Case 1: If $\alpha, \beta \leq \rho$ then $\pi_\alpha(p) = p$ and $\pi_\beta(q) = q$. So, using Lemma 5.17 and the assumption $\rho \leq \gamma$, we see that $f_\gamma(\pi_\alpha(p), \pi_\beta(q)) = f_\gamma(p, q) < 2\gamma$.

Case 2: If $\alpha \leq \rho < \beta$ then $\pi_\alpha(p) = p$ and Lemma 5.17 implies that $f_\beta(p, q) < 2\beta$.

$$\begin{aligned} f_\gamma(\pi_\alpha(p), \pi_\beta(q)) &= f_\gamma(p, \pi_\beta(q)) \\ &\leq f_\beta(p, \pi_\beta(q)) + 2(\gamma - \beta) \\ &\leq f_\beta(p, q) + 2(\gamma - \beta) \\ &< 2\gamma. \end{aligned}$$

Case 3: If $\rho < \alpha, \beta$ then Lemma 5.17 implies that $f_\alpha(p, q) < 2\alpha$.

$$\begin{aligned}
 f_\gamma(\pi_\alpha(p), \pi_\beta(q)) &\leq f_\beta(\pi_\alpha(p), \pi_\beta(q)) + 2(\gamma - \beta) \\
 &\leq f_\beta(\pi_\alpha(p), q) + 2(\gamma - \beta) \\
 &\leq f_\alpha(\pi_\alpha(p), q) + 2(\gamma - \beta) + 2(\beta - \alpha) \\
 &\leq f_\alpha(p, q) + 2(\gamma - \beta) + 2(\beta - \alpha) \\
 &< 2\gamma.
 \end{aligned}$$

■

Lemma 5.21 *Let (p, q) be an edge of R_δ with $\alpha, \beta \leq \frac{\delta}{1+\epsilon}$, then $(\pi_\alpha(p), \pi_\beta(q)) \in R_{\kappa\delta}$, where $\kappa = 1 + \frac{\sqrt{1+t^2}\epsilon}{1-\epsilon}$.*

Proof: First, note that the projection functions satisfy the following inequalities.

$$\begin{aligned}
 d_X(p, \pi_\alpha(p)) &\leq \epsilon\alpha \leq \frac{\epsilon\delta}{1-\epsilon} \\
 d_X(q, \pi_\beta(q)) &\leq \epsilon\beta \leq \frac{\epsilon\delta}{1-\epsilon}
 \end{aligned}$$

Applying the triangle inequality, the definition of an edge in R_δ , and Lemma 4.9, we get,

$$\begin{aligned}
 d_X(\pi_\alpha(p), \pi_\beta(q)) &\leq d_X(p, q) + \frac{2\epsilon\delta}{1-\epsilon} \\
 &< \left(r_p(\delta) + \frac{\epsilon\delta}{1-\epsilon}\right) + \left(r_q(\delta) + \frac{\epsilon\delta}{1-\epsilon}\right) \\
 &\leq r_{\pi_\alpha(p)}\left(\delta + \frac{\sqrt{1+t^2}\epsilon}{1-\epsilon}\delta\right) + r_{\pi_\beta(q)}\left(\delta + \frac{\sqrt{1+t^2}\epsilon}{1-\epsilon}\delta\right) \\
 &\leq r_{\pi_\alpha(p)}(\kappa\delta) + r_{\pi_\beta(q)}(\kappa\delta).
 \end{aligned}$$

This is precisely the condition sufficient to guarantee that $(\pi_\alpha(p), \pi_\beta(q)) \in R_{\kappa\delta}$ as desired. ■

The following two lemmas follow easily from repeated application of the preceding lemmas.

Lemma 5.22 *Two projections π_α and π_β induce contiguous simplicial maps from Q_ρ to Q_β whenever $\rho \leq \beta$ and there exists i so that $\frac{\lambda_{i+1}}{\epsilon(1-\epsilon)} \leq \alpha \leq \beta \leq \frac{\lambda_i}{\epsilon(1-\epsilon)}$.*

Proof: Fix $\rho \leq \beta$ and take (p, q) an edge from Q_ρ . Given that Q_ρ and Q_β are cliques complexes, we show that the simplex σ generated by $\{\pi_\alpha(p), \pi_\alpha(q), \pi_\beta(p), \pi_\beta(q)\}$ is in Q_β and apply Lemma 2.12. We only need to prove that all edges of σ belongs to Q_β .

We apply Lemma 5.20, while replacing γ by β and β by α . Thus we obtain $(\pi_\alpha(p), \pi_\alpha(q)) \in Q_\beta$. Let us repeat this operation with $\alpha = \beta = \gamma$ and we get $(\pi_\beta(p), \pi_\beta(q)) \in Q_\beta$. The last two edges are given by replacing γ by β and choosing correctly the role of p and q . ■

Lemma 5.23 *Two projections π_α and π_β induce contiguous simplicial maps from $R_\delta \rightarrow R_{\kappa\delta}$ for any $\delta \in \mathbb{R}$, where $\kappa = 1 + \frac{\sqrt{1+t^2}\epsilon}{1-\epsilon}$ whenever $\alpha, \beta \leq \frac{\delta}{1-\epsilon}$.*

Proof: The result is obtained by replacing Lemma 5.20 by Lemma 5.21 in the previous proof. ■

5.3.4 Sparse filtrations and power distance functions

We now show that $\{T_\alpha\}$ approximates $\{R_\alpha\}$ in terms of persistent homology by proving Theorem 5.14. To do this we demonstrate a multiplicative interleaving between these filtrations, where the interleaving constant is

$$\kappa = 1 + \frac{\sqrt{1+t^2}\epsilon}{1-\epsilon}.$$

Specifically, we show that for all $\alpha \geq 0$, the diagrams in Figure 5.5 commutes at the homology level.

$$\begin{array}{ccc} R_\alpha & \hookrightarrow & R_{\kappa\alpha} \\ \uparrow & \nearrow & \uparrow \\ T_\alpha & \hookrightarrow & T_{\kappa\alpha} \end{array} \qquad \begin{array}{ccc} R_\alpha & \hookrightarrow & R_{\kappa\alpha} \\ \uparrow & \searrow \pi_{\frac{\alpha}{1-\epsilon}} & \uparrow \\ T_\alpha & \hookrightarrow & T_{\kappa\alpha} \end{array}$$

Figure 5.5 – Diagrams with contiguous simplicial maps between $\{R_\alpha\}$ and $\{T_\alpha\}$

The diagrams on the left is commutative as the only maps involved are inclusions. For the diagram on the right, we first need to check that the projection $\pi_{\frac{\alpha}{1-\epsilon}}$ indeed induces a simplicial map from R_δ to $T_{\kappa\delta}$.

Lemma 5.24 *For all $\alpha > 0$, the projection $\pi_{\frac{\alpha}{1-\epsilon}}$ induces a simplicial map from $R_\alpha \rightarrow T_{\kappa\alpha}$, where $\kappa = 1 + \frac{\sqrt{1+t^2}\epsilon}{1-\epsilon}$.*

Proof: We show that for each edge $(p, q) \in R_\alpha$, there is a corresponding edge $(\pi_{\frac{\alpha}{1-\epsilon}}(p), \pi_{\frac{\alpha}{1-\epsilon}}(q)) \in R_{\kappa\alpha} \cap Q_{\frac{\alpha}{1-\epsilon}}$. Since the latter complex is a clique complex, this will imply that for all $\sigma \in R_\alpha$, we have $\pi_{\frac{\alpha}{1-\epsilon}}(\sigma) \in R_{\kappa\alpha} \cap Q_{\frac{\alpha}{1-\epsilon}} \subseteq T_{\kappa\alpha}$ as desired. First, $(\pi_{\frac{\alpha}{1-\epsilon}}(p), \pi_{\frac{\alpha}{1-\epsilon}}(q)) \in R_{\kappa\alpha}$ as a direct consequence of Lemma 5.23.

Next, we show that $(\pi_{\frac{\alpha}{1-\epsilon}}(p), \pi_{\frac{\alpha}{1-\epsilon}}(q)) \in Q_{\frac{\alpha}{1-\epsilon}}$. It suffices to prove that $f_{\frac{\alpha}{1-\epsilon}}(\pi_{\frac{\alpha}{1-\epsilon}}(p), \pi_{\frac{\alpha}{1-\epsilon}}(q)) < \frac{2\alpha}{1-\epsilon}$. We use Lemma 5.19 at the first line and the fact that $(p, q) \in R_\alpha$ for the fourth line.

$$\begin{aligned} f_{\frac{\alpha}{1-\epsilon}}(\pi_{\frac{\alpha}{1-\epsilon}}(p), \pi_{\frac{\alpha}{1-\epsilon}}(q)) &\leq f_{\frac{\alpha}{1-\epsilon}}(p, q) \\ &= d_{\mathbb{X}}(p, q) + s_p\left(\frac{\alpha}{1-\epsilon}\right) + s_q\left(\frac{\alpha}{1-\epsilon}\right) \\ &\leq d_{\mathbb{X}}(p, q) + \frac{2\epsilon\alpha}{1-\epsilon} \\ &< 2\alpha + \frac{2\epsilon\alpha}{1-\epsilon} \\ &= \frac{2\alpha}{1-\epsilon} \end{aligned}$$

■

Now, we give conditions for when two projections induce contiguous simplicial maps between the sparse weighted Rips complexes T_δ and $T_{\kappa\delta}$.

Lemma 5.25 *Two projections π_α and π_β induce contiguous simplicial maps from $T_\delta \rightarrow T_{\kappa\delta}$, where $\kappa = 1 + \frac{\sqrt{1+t^2}\epsilon}{1-\epsilon}$ whenever $\alpha, \beta \leq \frac{\delta}{1-\epsilon}$ and there exists i so that $\frac{\lambda_{i+1}}{\epsilon(1-\epsilon)} \leq \alpha \leq \beta \leq \frac{\lambda_i}{\epsilon(1-\epsilon)}$.*

Proof: We simply observe for any $\sigma \in T_\delta$, that $\sigma \in Q_\rho$ for some $\rho \leq \delta$. If $\rho \leq \beta$ then Lemma 5.22 implies $\pi_\alpha(\sigma) \cup \pi_\beta(\sigma) \in Q_\beta$. Otherwise $\pi_\alpha(\sigma) \cup \pi_\beta(\sigma) = \sigma \in Q_\rho$. So in either case, we have $\pi_\alpha(\sigma) \cup \pi_\beta(\sigma) \in S_{\kappa\delta}$. Now, by Lemma 5.23, we have $\pi_\alpha(\sigma) \cup \pi_\beta(\sigma) \in R_{\kappa\delta}$. So, we have $\pi_\alpha(\sigma) \cup \pi_\beta(\sigma) \in R_{\kappa\delta} \cap S_{\kappa\delta} = T_{\kappa\delta}$ as desired. ■

We can now give the proof of the interleaving which will imply the desired approximation of the persistent homology.

Lemma 5.26 *For all $\alpha > 0$, the diagrams in Figure 5.5 commutes at the homology level.*

Proof: By Lemma 5.23, the projection $\pi_{\frac{\alpha}{1-\epsilon}}$ and the inclusion π_0 are contiguous and thus produce identical homomorphisms at the homology level thanks to Theorem 2.11. For the lower triangle it will suffice to show that homomorphism induced by $\pi_{\frac{\alpha}{1-\epsilon}}$ commutes with the one produced by the inclusion π_0 . Let $\phi_i = \pi_{\frac{\lambda_i}{1-\epsilon}}$ for $i = 1, \dots, n+1$. Now, Lemma 5.25 implies that ϕ_i and ϕ_{i+1} are contiguous. So, choosing k such that $\lambda_k \leq \epsilon\alpha < \lambda_{k-1}$, we can apply Lemma 5.25 repeatedly to conclude that

$$\pi_{0*} = \phi_{n+1*} = \phi_{n*} = \dots = \phi_{k*} = \pi_{\frac{\alpha}{1-\epsilon}*}$$

where ϕ_{i*} is the homomorphism induced by ϕ_i . ■

The commutativity of these diagrams at the homology level is equivalent to the existence of a κ multiplicative interleaving between the two persistence modules. Theorem 5.14 is obtained by applying Corollary 2.25.

5.4 Experimental illustration

The algorithm using the distance to a measure and the sparse weighted Rips to compute persistence diagrams has been implemented. We used the ANN library [82] for the k -nearest neighbours search and code from Zomorodian following [99] for the persistence. The topology of the union of balls is acquired through the α -shapes implementation from the CGAL library [43].

We illustrate our results from three different perspectives: the quality of the approximation, the stability of the diagrams with respect to noise, and the size of the filtration after sparsification.

Datasets

For the first two parts, we consider the set of points in \mathbb{R}^3 obtained by sampling regularly the skeleton of the unit cube with 116 points. Then we add four noise points in the centre of four of its faces such that two opposite faces are empty. This is the example given in Section 1.2.

We would like to compute the persistence diagram of the skeleton of the cube. We write this diagram $\text{Dgm}(\text{Skel})$. It contains five homology generators in dimension 1 and one in dimension 2. Its barcode representation is given in Figure 5.7.

For sparsification, we use a slightly bigger dataset composed of 10000 points regularly distributed on a curve rolled around a torus. The point set is shown in Figure 5.8.

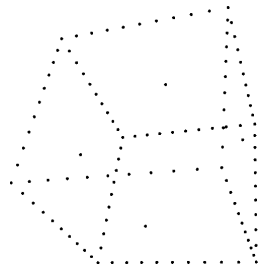


Figure 5.6 – Skeleton of a cube with outliers

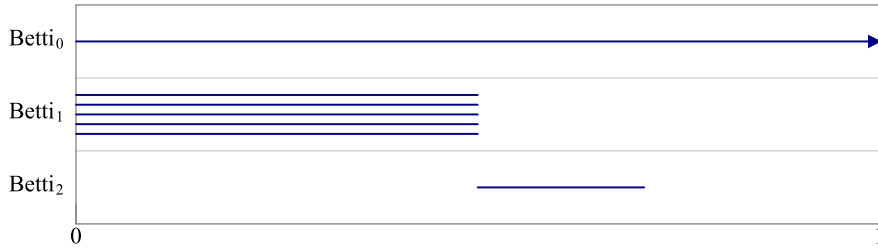


Figure 5.7 – Persistence diagram of a cube skeleton without outliers

Approximation We work on the cube skeleton with a mass parameter m such that $k = mn = 5$. We consider the empirical measure μ on the point cloud. The persistence diagram of $d_{\mu,m}$ is given in Figure 5.9.

The diagrams obtained with our various approximations have very similar looks. We only show the one obtained with the sparse Rips filtration with a parameter $\epsilon = 0.5$ in Figure 5.10. To compare diagrams, we use the bottleneck distances between the diagrams. Figure 5.11 shows the distance matrix between the various diagrams, while Figure 5.12 shows some bottleneck distances between persistence diagrams of different dimensions. Note that $\text{Dgm}(d_P)$ corresponds to the diagram obtained by using the distance function to the point cloud.

The largest difference is between $\text{Dgm}(\text{Skel})$ and $\text{Dgm}(d_{\mu,m})$. This is partly due to an effect of shifting while using the distance to a measure. After this initial shift, the distance are small compared to the theoretical bounds. Notice that the different steps of the approximation do not have the same effect on all dimensions.

All diagrams obtained by the different approximations are closer to $\text{Dgm}(\text{Skel})$ than the persistence diagram of the distance to the point cloud, $\text{Dgm}(d_P)$ given in Figure 5.13. For inference purposes, one crucial parameter is the *signal-to-noise ratio*. We define it as the ratio between the smallest lifespan of topological feature we aim to infer and the longest lifespan of noise features. A ratio of 1 corresponds to a signal that is not differentiable from the noise and ∞ corresponds to a noiseless diagram. In our example, only the dimensions 1 and 2 are relevant as the dimension 0 diagram corresponding to connected components has only one relevant feature and its lifespan is infinite. Results are listed in Figure 5.14.

Signal-to-noise ratios are clearly better than the one of $\text{Dgm}(d_P)$. Some of the approximation steps improve the ratio. This is due to two phenomena.

When one goes from $d_{\mu,m}$ to $d_{\mu,m}^P$, the filtration eliminates the cells of the k^{th} order Voronoi

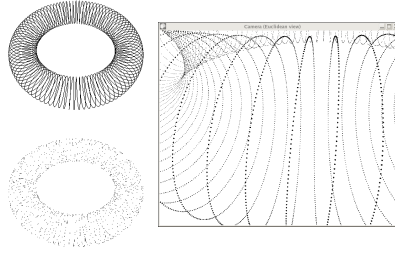


Figure 5.8 – Spiral on a torus

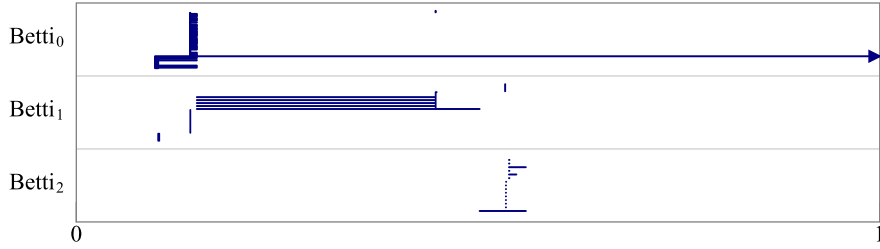


Figure 5.9 – $\text{Dgm}(d_{\mu,m})$ for the cube skeleton with outliers with $k = 5$

diagram that are far from the point cloud. These cells induce local minima that produce noise features in the diagrams. Removing them cleans parts of the diagram. The same phenomenon happens with the witnessed k -distance previously mentioned.

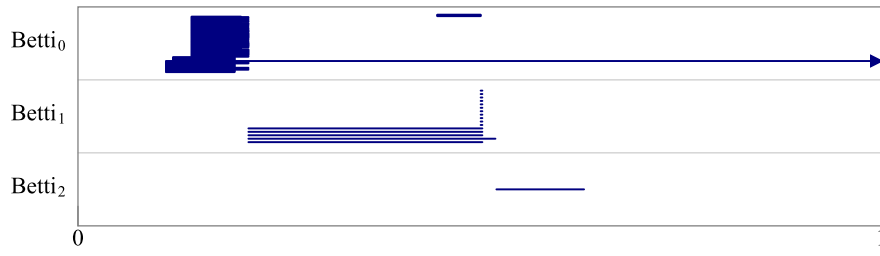
Using the Rips filtration instead of the Čech also reduces some noise. It eliminates artefacts from simplices that are introduced and almost immediately killed in the Čech complex due to balls that have pairwise intersections but no common intersection.

Stability

The weighted Rips filtration is stable with respect to noise. We illustrate this by studying the effect of an isotropic Gaussian noise on our skeleton of a cube. We consider three different standard deviations for our noise. Figure 5.15 shows the bottleneck distances between the persistence diagram of the sparse weighted Rips structure with the Gaussian noise and the one without Gaussian noise.

Unsurprisingly, the bottleneck distance is increasing with standard deviation of the noise. The signal-to-noise ratio shown in Figure 5.16 is more interesting.

Inferring correctly the homology of the cube skeleton is possible with standard deviation 0.05 and 0.1. Figure 5.17 shows the persistence diagram obtained with a standard deviation of 0.1. The ∞ in the 0.5 case in dimension 2 is not relevant as there is no noise but the feature is too small compared to the rest of the diagram as shown in Figure 5.18. Note that 0.5 corresponds to half the side of the cube, and thus, it is logical to be unable to retrieve any useful information. Some structure appears even with standard deviation as large as 0.5. The three bigger features in dimension 1 are relevant. However, we miss two elements and it is difficult to decide where to draw the frontier between relevant and irrelevant features.

Figure 5.10 – $\text{Dgm}(\{T_\alpha\})$ for the cube skeleton with outliers with $k = 5$ and $\epsilon = .5$

	$\text{Dgm}(Skel)$	$\text{Dgm}(d_{\mu,m})$	$\text{Dgm}(d_{\mu,m}^P)$	$\text{Dgm}(R_\alpha)$	$\text{Dgm}(T_\alpha)$	$\text{Dgm}(d_P)$
$\text{Dgm}(Skel)$	0	.1528	.1473	.1473	.1817	.25
$\text{Dgm}(d_{\mu,m})$.1528	0	.09872	.0865	.1183	.2543
$\text{Dgm}(d_{\mu,m}^P)$.1473	.09872	0	.0459	.1084	.2642
$\text{Dgm}(R_\alpha)$.1473	.0865	.0459	0	.1128	.2598
$\text{Dgm}(T_\alpha)$.1817	.1183	.1084	.1128	0	.2484
$\text{Dgm}(d_P)$.25	.2543	.2642	.2598	.2484	0

Figure 5.11 – Matrix of distances for the bottleneck distance

Sparsification efficiency

We introduced sparsification in Section 5.3.4 to reduce the size of the Rips filtration. The method introduced a new parameter ϵ , and the size of the filtration depends heavily on ϵ . The evolution of the size of the filtration depending on the parameter ϵ is given in Figure 5.19 for the sampling of the spiral.

The minimum size is reached around $\epsilon = .83$. This minimum depends on the structure of the dataset. For example, considering a set of points uniformly sampled in a square, we obtain a size that is monotonic.

The filtration size is nearly constant after a rapid decrease. In this example, the size is of order 10^7 simplices for an input of 10^5 vertices. Computing persistent homology is tractable for any value in this range.

Finally, we illustrate the dependence on the intrinsic dimension. Figure 5.20 shows the number of simplices in the sparse Weighted Rips filtration for different values of ϵ and number of sample points in \mathbb{R}^2 . On the left, the points are sampled uniformly in a unit square and hence have intrinsic dimension 2. On the right, the points are sampled uniformly on a unit circle and thus have intrinsic dimension 1. The lower intrinsic dimension of the second example explains the difference of size between the two sparsified filtrations.

Dgm(A)	Dgm(B)	dim 0	dim 1	dim 2
Dgm($Skel$)	Dgm($d_{\mu,m}$)	.05202	.1528	.1495
Dgm($d_{\mu,m}$)	Dgm($d_{\mu,m}^P$)	.09872	.0195	.0972
Dgm($d_{\mu,m}^P$)	Dgm($R_\alpha(P, d_{\mu,m})$)	.0007	.0044	.0459
Dgm($R_\alpha(P, d_{\mu,m})$)	Dgm($T_\alpha(P, d_{\mu,m})$)	.0872	.1128	.0026
Dgm($Skel$)	Dgm($d_{\mu,m}^P$)	.0405	.1473	.0982
Dgm($Skel$)	Dgm($T_\alpha(P, d_{\mu,m})$)	.1026	.1817	.098
Dgm($Skel$)	Dgm(d_P)	.25	.2071	.1481

Figure 5.12 – Bottleneck distances between diagrams

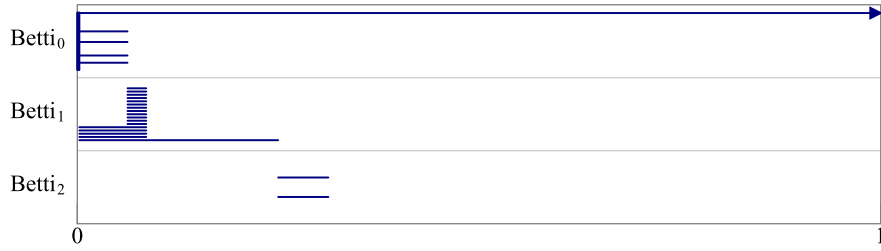


Figure 5.13 – Dgm(d_P) for the cube skeleton with outliers

Diagram	dim 1	dim 2
Dgm($Skel$)	∞	∞
Dgm($d_{\mu,m}$)	247	2.74
Dgm($d_{\mu,m}^P$)	69.8	43
Dgm($R_\alpha(P, d_{\mu,m})$)	∞	∞
Dgm($T_\alpha(P, d_{\mu,m})$)	132	∞
Dgm(d_P)	5.66	1

Figure 5.14 – Signal to noise ratios

Standard deviation	.05	.1	.5
d_b in dimension 1	.1469	.2261	.2722
d_b in dimension 2	.047	.0914	.1046

Figure 5.15 – d_b between Dgm($\{T_\alpha\}$) with and without Gaussian noise

Standard deviation	0	.05	.1	.5
Ratio in dimension 1	132	8.27	3.17	1.04
Ratio in dimension 2	∞	∞	100.2	∞

Figure 5.16 – Signal to noise ratio of Dgm($\{T_\alpha\}$) depending on noise intensity

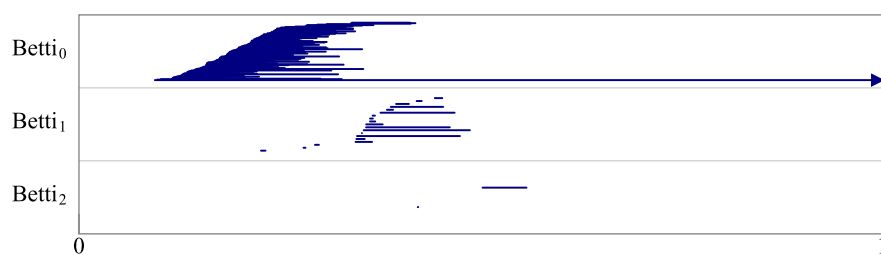


Figure 5.17 – Persistence diagram of $\{T_\alpha\}$ with $k = 5$, $\epsilon = 0.5$ and a Gaussian noise with standard deviation 0.1

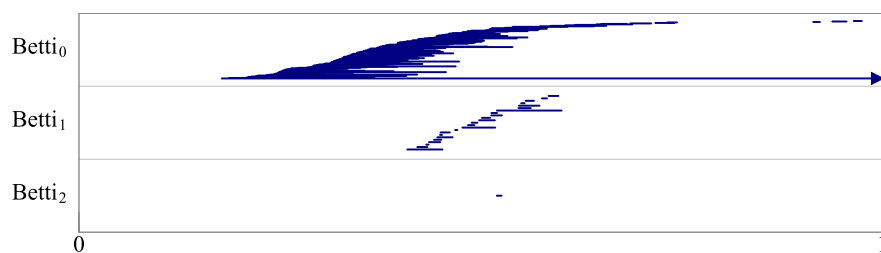


Figure 5.18 – Persistence diagram of $\{T_\alpha\}$ with $k = 5$, $\epsilon = .5$ and a Gaussian noise with standard deviation .5

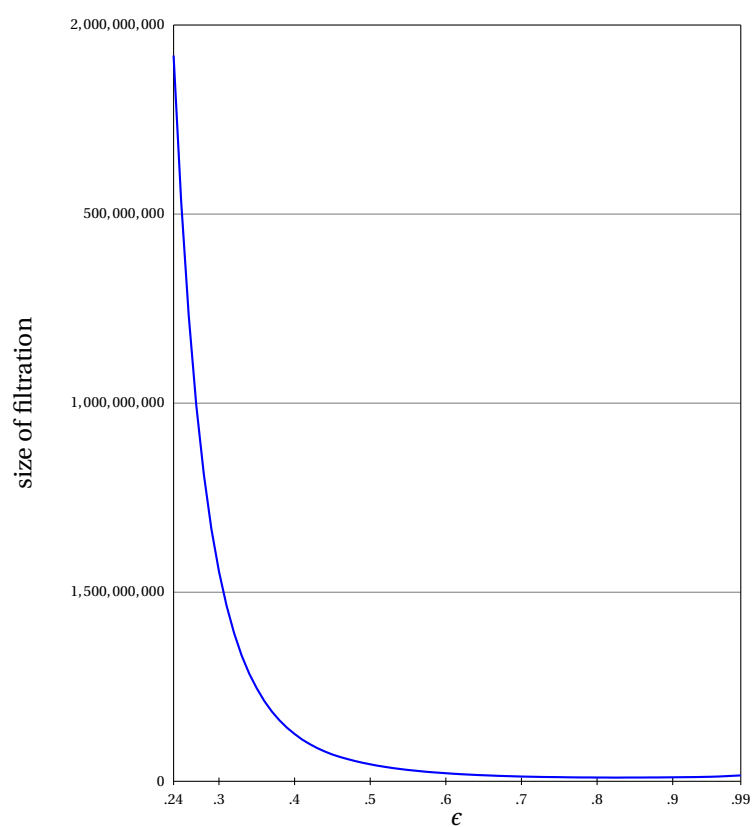
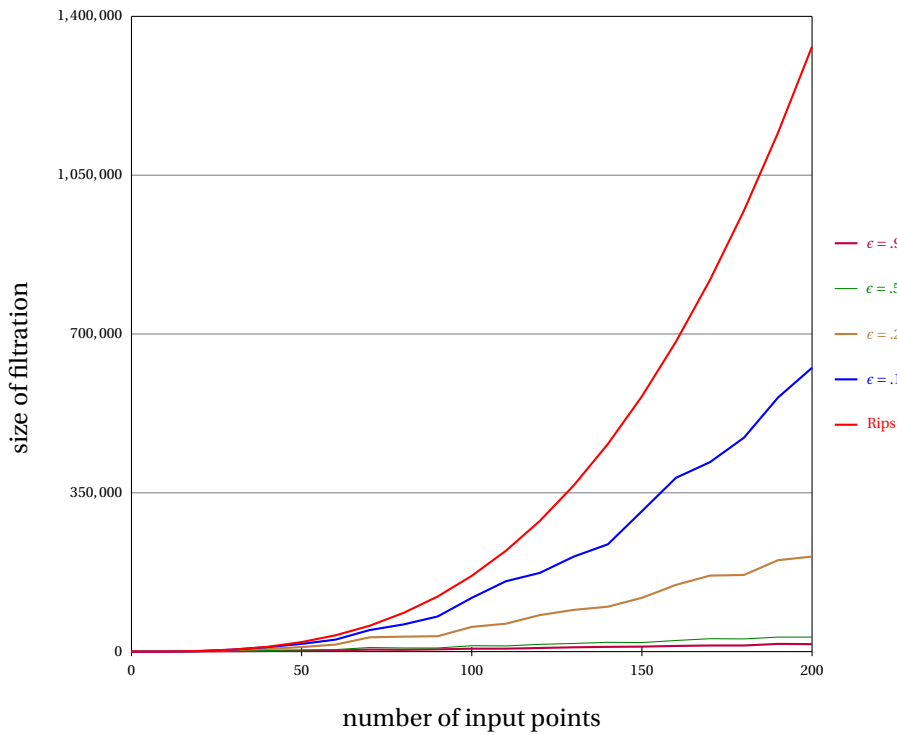
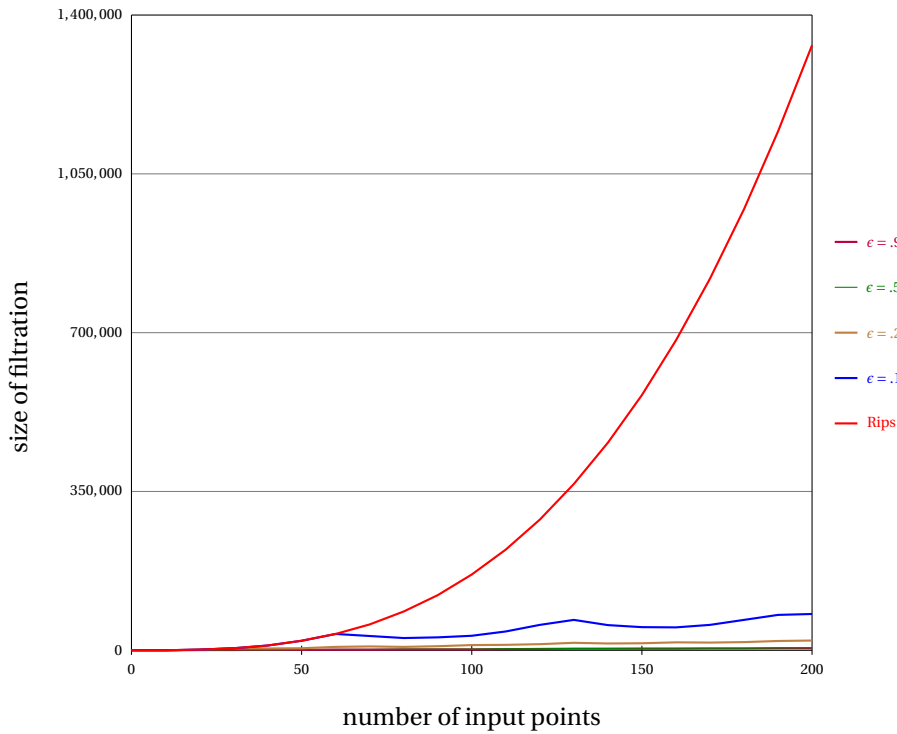


Figure 5.19 – Size of the filtration depending on ϵ for the spiral



(a) Points sampled in a square



(b) Points sampled on a circle

Figure 5.20 – Size of the filtration depending on the number of input points

6 Specific noise conditions for the distance to a measure

Most results in topological data analysis are obtained under the condition that two sets are close with respect to the Hausdorff distance or that two measures are close with respect to the Wasserstein distance. Previous results for the distance to a measure relied on the Wasserstein distance. However, this distance is not always tight to indicate the proximity of distance to measure functions. In this chapter, we discuss the interest of using Wasserstein distances and their limitations. Then, we introduce new set of noise conditions and motivate them with the example of clutter noise. Finally, we show how these new conditions relate to other classical conditions and translates into guarantees on persistence diagrams.

6.1 Counter examples for the Wasserstein noise condition

Stability results for the distance to a measure were originally stated using the Wasserstein distance. This is well-suited to objects that are probability measures. However, the classical setting for topological data analysis is not built upon a measure. We usually consider an object M as the *ground truth*. In most theoretical results, this object is assumed to be a compact set or a Riemannian submanifold.

The object M represents what we want to find at the end of our analysis. For example, this can be a human body. We scan it and obtain a point cloud which is a representation of this body, with the possible presence of noise and defects in the scan, and we want to reconstruct the shape of the body. We aim to reconstruct a shape having the same topology as the human body and being as close as possible to M , for example in the Hausdorff distance.

The simplest condition on a sampling to obtain good results in reconstruction is that the input point cloud P lies close to the ground truth M in the Hausdorff distance. This distance is very sensitive to outliers. Only one misplaced point is necessary to have a large Hausdorff distance. The Wasserstein distance aims to take into account that some points can be bad, but it presents two major difficulties.

First, the Wasserstein distance is a distance between probability measures. The setting is different and we have to adapt assumptions. A natural way to handle the point cloud P is to take the empirical measure on it. For M , it is more difficult, especially if we do not have a generative model for P . If M is a compact Riemannian submanifold of R^d , we can for example

take the uniform measure on M . However, if the sampling is strongly uneven, we could end up with a large Wasserstein distance, while the topological analysis algorithm would work well and the Hausdorff distance would be small.

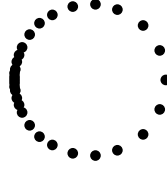


Figure 6.1 – Uneven sampling of a circle

Figure 6.1 shows a circle with an uneven sample. The ground truth is the circle M and we only have access to the point cloud P . In this case, the Hausdorff distance between P and M is small and given by half the largest gap ϵ between two consecutive points on the circle. The Wasserstein distance between the empirical measure on P and the uniform measure on M is much bigger due to the difference of sampling density, and it can be arbitrarily close to $2\sqrt{2\pi}r$, where r is the radius of the circle. At the limit, all the mass is at one point on the circle and the Wasserstein distance between the Dirac $\delta_{-1,0}$ and the uniform measure on the circle μ_C is given by:

$$W_2(\delta_{(-1,0)}, \mu_C) = 2\sqrt{r^2 \int_0^\pi \sin^2(\theta) + (1 + \cos(\theta))^2 d\theta} = 2\sqrt{\int_0^\pi 2r^2(1 + \cos(\theta)) d\theta} = 2\sqrt{2\pi}r.$$

The second difficulty is that two measures can be arbitrarily far from each other in the Wasserstein distance, while the relative difference between distances to them is arbitrarily small. Let us consider the metric space \mathbb{R} and the Dirac measure δ_0 at the origin. Now, consider the sequence of measures $\nu_n = \frac{n-1}{n}\delta_0 + \frac{1}{n}\delta_{\frac{3}{n^2}}$ for $n \geq 2$. This means that we take a small mass and send it further and further away, more quickly than the moved mass decreases. If the mass considered for the distance to the measures is positive and fixed, then the relative error between $d_{\delta_0, m}$ and $d_{\nu_n, m}$ tends to 0.

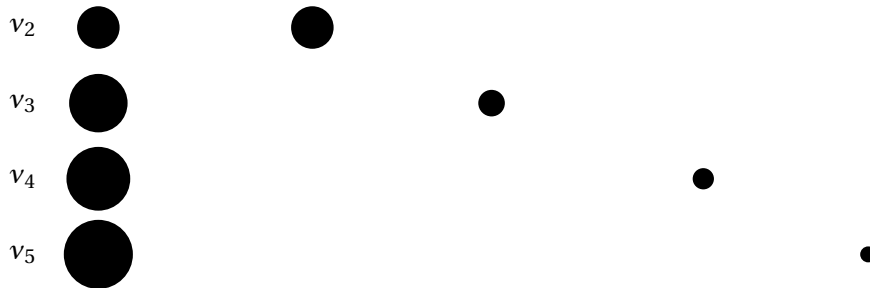


Figure 6.2 – First ν_n for $n \geq 2$, where the size of a circle corresponds to the mass at its centre

There exists an unique transport plan between δ_0 and ν_n . The Wasserstein distance has value $W_2(\delta_0, \nu_n) = n$. Consider a mass $m \in]0, 1[$ and an integer n such that $m > \frac{1}{n}$. Then, for any

$x \leq \frac{n^{\frac{3}{2}}}{2}$, $d_{\delta_0,m}(x) = d_{v_n,m}(x) = |x|$ and for $x \geq \frac{n^{\frac{3}{2}}}{2}$, $d_{\delta_0,m}(x) = x$ and $d_{v_n,m}(x) = \sqrt{\frac{1}{m} \left(\left(m - \frac{1}{n}\right) x^2 + \frac{1}{n} (n^{\frac{3}{2}} - x)^2 \right)}$. Hence:

$$\begin{aligned}
 \frac{d_{\delta_0,m}(x) - d_{v_n,m}(x)}{d_{\delta_0,m}(x)} &= \frac{x - \sqrt{x^2 - \frac{2n^{\frac{1}{2}}}{m}x + \frac{n^2}{m}}}{x} \\
 &\leq 1 - \sqrt{1 - \left(\frac{4}{mn} + \frac{2}{mn}\right)} \\
 &= \frac{3}{mn} + o\left(\frac{1}{n}\right)
 \end{aligned}$$

Thus, when n tends to ∞ , the relative error between the functions tends to 0. This is due to the fact that the Wasserstein distance is only worst-case tight for distances to measures. The relation $\|d_{\mu,m} - d_{v,m}\|_{\infty} \leq \frac{1}{\sqrt{m}} W_2(\mu, \nu)$ is an equality only in some worst-case example contrarily to $\|d_P, d_K\|_{\infty} = d_H(P, K)$ which is always tight as shown in Lemma 2.29.

In the presence of background noise, the Wasserstein distance and current stability results present another drawback. The sub-level sets of the distance to a measure can be trivial topological balls due to the background noise before recovering all features.

Consider a mixture of two measures. For the first one, consider the uniform measure μ_M on eight circles of radius 1 regularly centred on a circle of radius 10. The support of μ_M is a manifold M and constitutes the ground truth. However, we have a uniform noise materialised by a measure μ_B uniform on the square of side 1000. We only have access to $\mu = \lambda\mu_M + (1-\lambda)\mu_B$ for some $\lambda \in [0, 1]$.

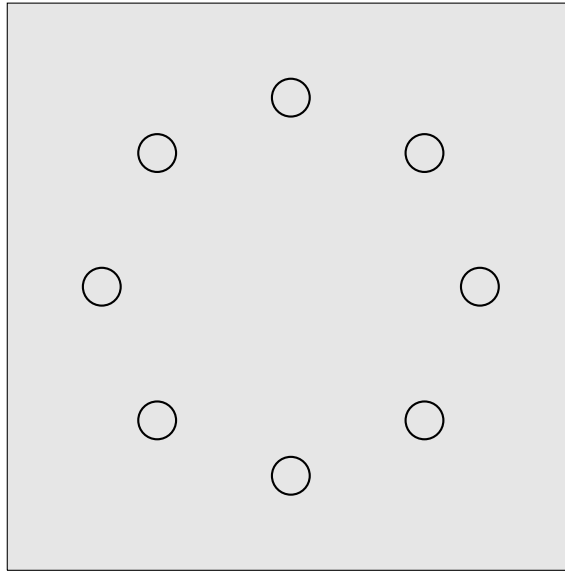


Figure 6.3 – Example with eight circles, the bounding box being reduced

M has two structures in dimension 1. First, there is the set of eight circles. At larger scales, it

becomes the larger circle. However, due to the background noise μ_B , we are unable to recover anything using the distance to the support of the mixed measure μ . Using the distance to μ with a small mass parameter m and if λ is sufficiently large, we can build a persistence diagram containing features for every small circle but not for the large one. Figure 6.4 shows the value of $d_{\mu,m}$ along a line crossing the centres o_1 and o_2 of two neighbouring small circles. Example of values for λ and m are given in Section 6.3.1.

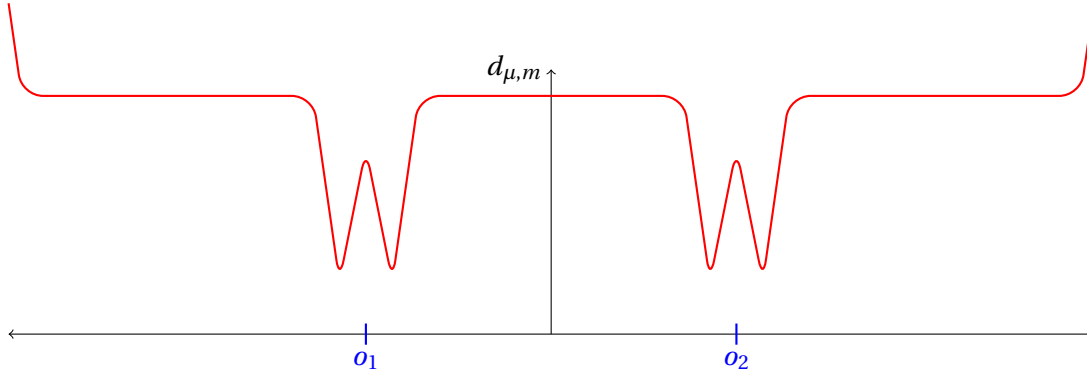


Figure 6.4 – Value $d_{\mu,m}$ along a line between two neighbouring centres

Using the distance to a measure, we are unable to recover a good approximation of the whole persistence diagram. The sub-level sets of $d_{\mu,m}$ have the topology type of the eight small circles at the beginning, but as we reach the value where $d_{\mu,m}$ is constant, the topology of the sub-level sets becomes the topology of a ball. The large cycle does not appear but we can recover the part of the persistence diagram corresponding to low scales, id est, small values of the filtration parameter. Our new noise conditions aim at guaranteeing such an approximation.

6.2 New noise conditions for distances to measures

We introduce a new set of noise conditions to correct the above flaws. Our noise conditions consider the point set as a measure and relates it to the ground truth M using two parameters and distance to measure functions.

Definition 6.1 Let $M \subset \mathbb{R}^d$ be a manifold and let μ be a probability measure. For a fixed $m \in [0, 1]$, μ is an (ϵ, r) -sample of M if:

$$\epsilon \geq \sup_{x \in M} d_{\mu,m}(x)$$

$$r \leq \sup\{\ell \in \mathbb{R} \mid \forall x \in \mathbb{R}^d, d_{\mu,m}(x) < \ell \implies d(x, M) \leq d_{\mu,m}(x) + \epsilon\}$$

By extension, if $P \subset \mathbb{R}^d$ is a point set, we say that P is an (ϵ, r) -sample of M if the empirical measure on P is an (ϵ, r) -sample of M .

The parameter ε captures the distance to the empirical measure for points in M and thus tells us how dense the sampling measure is around the manifold M . The parameter r tells us how far away we can go from the manifold without having the noise mask relevant information. Remark that if a point set is an (ε, r) -sample of M then it is an (ε', r') -sample for any $\varepsilon' \geq \varepsilon$ and $r' \leq r$. Be careful that this (ε, r) -sampling condition is different from the classical (ε, δ) condition [37] where δ is a packing condition.

For convenience, denote the distance function to the manifold M by $d_M : \mathbb{R}^n \rightarrow \mathbb{R}, x \mapsto d(x, M)$. We have the following interleaving between the sub-level sets of $d_{\mu, m}$ and d_M .

Lemma 6.2 *Let μ be an (ε, r) -sample of a manifold M . Then,*

$$\forall \alpha < r - \varepsilon, d_M^{-1}([-\infty, \alpha]) \subset d_{\mu, m}^{-1}([-\infty, \alpha + \varepsilon]) \subset d_M^{-1}([-\infty, \alpha + 2\varepsilon]).$$

Proof: Let x be a point such that $d(x, M) \leq \alpha$. There exists $y \in M$ such that $d(y, x) \leq \alpha$. By definition of ε , $d_{\mu, m}(y) \leq \varepsilon$. Given that the distance to a measure is a 1-Lipschitz function, $d_{\mu, m}(x) \leq \varepsilon + \alpha$.

Let now x be a point such that $d_{\mu, m}(x) \leq \alpha + \varepsilon \leq r$. By definition of r , we have $d(x, M) \leq d_{\mu, m}(x) + \varepsilon \leq \alpha + 2\varepsilon$ which concludes the argument. ■

This provides a partial interleaving between the sub-level sets of the distance to μ and the offsets of the manifold M . Observe that this relation is very similar to the one obtained in lemma 2.29 when two compact sets A and B have Hausdorff distance of at most ε :

$$\forall \alpha, d_A^{-1}([-\infty, \alpha]) \subset d_B^{-1}([-\infty, \alpha + \varepsilon]) \subset d_A^{-1}([-\infty, \alpha + 2\varepsilon])$$

Various results on inference are given for point sets that are close to the underlying manifold in the Hausdorff distance. An algorithm, using the distance to P and requiring $d_H(P, M) \leq \varepsilon$, has the same guarantee using the distance to a measure and an (ε, r) -sample of M as long as $r > 2\varepsilon$. It gives a set of reconstruction methods where we can use existing algorithms, for example [11, 13, 34, 35, 71], and replace the distance to the input point cloud by the distance to an empirical measure.

6.3 Relation to other noise conditions

To show the interest of the distance to a measure and of our noise conditions, we study the relation between other noise conditions and these new ones. We remark that our noise conditions encompass many other existing noise conditions. While the parameter ε is natural, the parameter r may appear to be artificial. It bounds the distances at which we can observe the manifold through the lens of the distance to a measure to obtain guarantees on persistence diagram inference. In most classical noise conditions, r is equal to ∞ and thus we obtain exactly the same relation as for the classical Hausdorff condition between the input point cloud P and the manifold M . One notable noise model where $r < \infty$ is when there is a uniform background noise in the ambient space \mathbb{R}^d , sometimes called *clutter noise*. In this case, r will depend on the density ratio between noise and information.

6.3.1 Clutter noise

A strong background noise can completely mess up a data analysis. Sometimes called *clutter noise*, it corresponds to an uniform background density over the area where we take measurements. We consider a restricted case to have an easier development.

We consider a compact Riemannian submanifold $M \subset \mathbb{R}^d$ of dimension d' . Let μ_M be the uniform probability measure on M and μ_B be the uniform probability measure on an Euclidean ball of radius R , such that the ball largely encloses M . We mean that the centre of mass of M and the centre of the ball are the same, and that R is much greater than the diameter of M . The measure we observe is $\mu = \lambda\mu_M + (1 - \lambda)\mu_B$, where $\lambda \in [0, 1]$.

Theorem 6.3 *Given a compact Riemannian submanifold $M \subset \mathbb{R}^d$ with positive reach r_M , positive strong convexity radius $\rho(M)$, and whose curvature is bounded by c_M and a mass parameter m sufficiently small, there exists $\epsilon \geq 0$ and $r > 0$, such that $\mu = \lambda\mu_M + (1 - \lambda)\mu_B$ is an (ϵ, r) -sample of M for any $\lambda \in]0, 1]$.*

Proof: Consider a ball \mathcal{B} of radius ρ enclosed in the large ball of radius R and that does not intersect M . Then, there exists C such that $\mu(\mathcal{B}) = (1 - \lambda)\mu_B(\mathcal{B}) = C\rho^d$. Let x be a point of M . Using Lemma 2.37, for any $a \leq \min(\rho(M); \frac{\pi}{\sqrt{c_M}})$, $\mu(B(x, a)) \geq \lambda\mathcal{C}_{d'}^{c_M} a^{d'} + Ca^d$. Fix $b = \min(\rho(M); \frac{\pi}{\sqrt{c_M}})$ and $m_0 = \lambda\mathcal{C}_{d'}^{c_M} b^{d'} + Cb^d$. For any $m \leq m_0$, we can thus bound $\delta_{\mu, m}(x)$ and therefore $d_{\mu, m}(x)$. Remark that, when $m \rightarrow 0$, $d_{\mu, m}(x) \rightarrow 0$. We fix m such that $m \leq m_0$ and $d_{\mu, m}(x) < r_M$ and we introduce $\epsilon = \sup_{x \in M} d_{\mu, m}(x)$.

Now consider a point x such that $d(x, M) = (\frac{m}{C})^{\frac{1}{d}}$. Then $d_{\mu, m}(x) = \sqrt{\frac{m^{1+\frac{2}{d}}}{(1+\frac{2}{d})C^{\frac{2}{d}}}} = r < r_M$. Let

y be a point such that $d_{\mu, m}(y) < r$. Therefore, $d(y, M) < (\frac{m}{C})^{\frac{1}{d}}$ and there exists x such that $d(x, M) = (\frac{m}{C})^{\frac{1}{d}} = d(y, M) + d(x, y)$. Hence $d_{\mu, m}(y) \geq r - d(x, y) = r - d(x, M) + d(y, M)$. And $d(y, M) \leq d_{\mu, m}(y) + (\frac{m}{C})^{\frac{1}{d}} - r$.

Fixing, $\epsilon = \max(\sup_{x \in M} d_{\mu, m}(x); (\frac{m}{C})^{\frac{1}{d}} - r)$, μ is an (ϵ, r) -sample of M . ■

Remark that the proof of existence is not enough for inference purposes. We need to have $r > 2\epsilon$ to obtain something interesting in Lemma 6.2. It is necessary to handle the bound case by case. We present the computation of interesting bounds for the example of Figure 6.3.

Fix $\lambda = 8 \times 10^{-4}$. Then each small circle has mass 10^{-4} and the measure μ_B has density 9.9992×10^{-7} on its support. We use $m = 5 \times 10^{-6}$. Let x be a point on one of the small circles. For $r \leq 1$,

$$\mu(B(x, a)) \geq \lambda\mu_M(B(x, a)) = \frac{a}{\pi} \sqrt{1 - \frac{a^2}{4}} \times 10^{-4} \geq \frac{a}{\pi} \times 10^{-4}.$$

Hence $\delta_{\mu, l}(x) \leq \pi l \times 10^4$ for $l \leq m$.

$$d_{\mu, m}(x) \leq \sqrt{\frac{1}{m} \int_0^m \pi^2 \times 10^8 l^2 dl} = \frac{5\pi \times 10^{-2}}{\sqrt{3}} \leq 9 \times 10^{-2}$$

Thus $\epsilon = \frac{5\pi \times 10^{-2}}{\sqrt{3}}$ works for our noise model. Let $y \in \mathbb{R}^2$ such that $\delta_{\mu, m}(y) \leq d_M(y)$. It implies that there exists a ball around y of radius a such that $\mu(B(x, a)) \geq m$ and $a \leq d_M(y)$. The

parameter r is greater than the infimum of the possible a for all y . For $a \leq d_M(y)$, $\mu(B(x, a)) = \pi a^2 9.9992 \times 10^{-7}$, which implies

$$a \geq \sqrt{\frac{m 10^7}{\pi}} = \sqrt{\frac{.5}{\pi}} \geq .39.$$

Moreover, such a point y always exist at distance $\sqrt{\frac{.5}{\pi}}$ of M and thus the parameter r is bounded:

$$.39 \leq r \leq \sqrt{\frac{.5}{\pi}} + \epsilon \leq .5.$$

6.3.2 Wasserstein noise condition

In the Wasserstein noise condition, the empirical measure μ on P is assumed to be close to the uniform measure μ_M on a Riemannian d' -manifold M in the Wasserstein distance. Assume that the curvature of M is bounded by c_M and that M has a positive strong convexity radius $\rho(M)$. Let $b = \min(\frac{\pi}{\sqrt{c_M}}; \rho(M))$. V_M denotes the volume of M while $\mathcal{C}_{d'}^{c_M} = \frac{4}{d'} \Gamma(\frac{1}{2})^{d'} \Gamma(\frac{d'}{2})^{-1} \left(\frac{\sqrt{c_M}}{\pi}\right)^{d'-1}$ is a constant. See Lemma 2.37 for details on $\mathcal{C}_{d'}^{c_M}$.

Theorem 6.4 *A measure μ with $W_2(\mu, \mu_M) \leq \sigma$ is an (ϵ, r) -sample of M for $m \leq \frac{\mathcal{C}_{d'}^{c_M} b^{d'}}{V_M}$,*

$$\epsilon \geq \frac{1}{\sqrt{1 + \frac{2}{d'}}} \left(\frac{V_M m}{\mathcal{C}_{d'}^{c_M}} \right)^{\frac{1}{d'}} + \frac{\sigma}{\sqrt{m}}, \text{ and } r = \infty.$$

Proof: Let μ_M be the uniform measure on M and μ a measure such that $W_2(\mu, \mu_M) \leq \sigma$. Using Theorem 3.15, we get $\|d_{\mu, m} - d_{\mu_M, m}\|_\infty \leq \frac{\sigma}{\sqrt{m}}$. Let us consider a point $x \in M$ and the ball $B(x, a)$ centred in x and of radius a . By definition of μ_M , for any $a \leq b$, using Lemma 2.37,

$$\mu_M(B(x, a)) = \frac{\mathcal{V}(x, a)}{V_M} \geq \frac{\mathcal{C}_{d'}^{c_M} a^{d'}}{V_M}$$

The pseudo-distance $\delta_{\mu_M, m}(x)$ can be bounded as long as $m \leq \frac{\mathcal{C}_{d'}^{c_M} b^{d'}}{V_M}$.

$$\delta_{\mu_M, m}(x) \leq a \leq \frac{\pi}{c_M} \leq \left(\frac{m V_M}{\mathcal{C}_{d'}^{c_M}} \right)^{\frac{1}{d'}}$$

And the relation propagates to the distance to μ_M .

$$\begin{aligned} d_{\mu_M, m}(x) &\leq \frac{1}{\sqrt{m}} \sqrt{\int_0^m \left(\frac{V_M}{\mathcal{C}_{d'}^{c_M}} l \right)^{\frac{2}{d'}} dl} \\ &\leq \frac{1}{\sqrt{1 + \frac{2}{d'}}} \left(\frac{V_M m}{\mathcal{C}_{d'}^{c_M}} \right)^{\frac{1}{d'}} \end{aligned}$$

Thus:

$$d_{\mu,m}(x) \leq \frac{1}{\sqrt{1 + \frac{2}{d'}}} \left(\frac{V_M m}{\mathcal{C}_{d'}^{c_M}} \right)^{\frac{1}{d'}} + \frac{\sigma}{\sqrt{m}}$$

The first part of our noise conditions is therefore verified for any $\epsilon \geq \frac{1}{\sqrt{1 + \frac{2}{d'}}} \left(\frac{V_M m}{\mathcal{C}_{d'}^{c_M}} \right)^{\frac{1}{d'}} + \frac{\sigma}{\sqrt{m}}$.

Moreover, for any $x \in \mathbb{R}^d$, $d_{\mu_M,m}(x) \geq d(x, M)$ because M is the support of μ_M . Thus:

$$d(x, M) \leq d_{\mu_M,m}(x) \leq d_{\mu,m}(x) + \frac{\sigma}{\sqrt{m}} \leq d_{\mu,m}(x) + \epsilon$$

and we have proved that μ is an (ϵ, ∞) -sampling of M . ■

Remark that the minimal value for ϵ is not 0 when σ is 0. To reach 0, we also need to decrease m in order to make the radius of the smallest ball containing a mass m tend to 0 at any point $x \in M$. From this bound for the Wasserstein noise condition, we derive an immediate result for the Gaussian noise model. Let us write $\mathcal{N}(0, \sigma^2)$ for the normal law on \mathbb{R}^d with zero mean and standard deviation σ and \star for the convolution operator.

Corollary 6.5 *The measure $\mu = \mu_M \star \mathcal{N}(0, \sigma^2)$ is an (ϵ, r) -sample of M for $m \leq \frac{\mathcal{C}_{d'}^{c_M} b^{d'}}{V_M}$,*

$$\epsilon \geq \frac{1}{\sqrt{1 + \frac{2}{d'}}} \left(\frac{V_M m}{\mathcal{C}_{d'}^{c_M}} \right)^{\frac{1}{d'}} + \frac{\sigma}{\sqrt{m}}, \quad \text{and } r = \infty.$$

Proof: It suffices to show that $W_2(\mu, \mu_M) \leq \sigma$, and then to apply Theorem 6.4. Denoting X the probability law corresponding to μ_M and Y , the law corresponding to $\mathcal{N}(0, \sigma^2)$, the probability measure μ corresponds to the law $X + Y$. Following [94, Definition 6.1],

$$W_2(\mu, \mu_M) = \inf \{ \sqrt{\mathbb{E}[|X - Z|^2]} \mid \text{law}(X) = \mu, \text{law}(Y) = \mu_M \} \leq \sqrt{\mathbb{E}[|X - X - Y|^2]} \leq \sigma.$$

■

6.3.3 Sampling by empirical measures

Data is usually given by a set P of n points sampled according to a measure μ . As we have seen previously, we consider the empirical measure μ_n on P . This measure has convergence properties to μ in Wasserstein metrics. Intuitively, if we sample more and more points according to μ , the empirical measure μ_n will be more and more similar to μ .

More formally, let us consider a probability measure μ on \mathbb{R}^d . Let X_1, X_2, \dots, X_n be independent identically distributed variables with probability law μ . Let $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ be their empirical measure. μ_n converges to μ in the second order Wasserstein metric [73, §1.1]:

Theorem 6.6 *Suppose that $c = \int |u|^{d+5} \mu(du) < \infty$, then there is a constant C depending only on c and the dimension d such that:*

$$\mathbb{E}[W_2(\mu, \mu_n)^2] \leq \frac{C}{n^{\frac{2}{d+4}}}.$$

This result can be combined with Theorem 6.4.

Corollary 6.7 *Let P_n be a set of n points obtained by an independent identically distributed sampling according to the uniform measure μ_M on a compact Riemannian submanifold $M \subset \mathbb{R}^d$ with curvature bounded by c_M and positive strong convexity radius $\rho(M)$. For any mass $m \leq \frac{c_M^{c_M} b^{d'}}{V_M}$ where $b = \min(\frac{\pi}{\sqrt{c_M}}; \rho(M))$, there exists a constant C' depending only on M and the dimension d such that, with probability at least $\frac{1}{2}$, μ_n , the empirical measure on P_n , is an (ϵ, r) -sample of M for:*

$$\epsilon \geq \frac{1}{\sqrt{1 + \frac{2}{d'}}} \left(\frac{V_M m}{c_M^{c_M}} \right)^{\frac{1}{d'}} + \frac{C'}{n^{\frac{1}{d+4}} \sqrt{m}}, \quad \text{and } r = \infty.$$

Proof: Theorem 6.6 implies that, with probability at least $\frac{1}{2}$, $W_2(\mu_n, \mu_M) \leq \frac{\sqrt{2C}}{n^{\frac{1}{d+4}}}$. We apply Theorem 6.4 and write $C' = \sqrt{2C}$. ■

This remark justifies the use of Wasserstein distances when working with empirical measures. Convergence rates for other orders of the Wasserstein metric exists. We refer the interested reader to [53, 55, 73] for more details on the behaviour of empirical measures.

6.3.4 Discrete results for Hausdorff noise condition

Given a compact Riemannian manifold M , the set of points P is assumed to be a ρ -sampling of M . This means that for each point of M , there exists a point of P that is at distance at most ρ . Assuming that the sampling is lying on the manifold and conditions are given in the Riemannian metric, we obtain:

Theorem 6.8 *Let M be a connected compact Riemannian manifold embedded in \mathbb{R}^d and let P be a Riemannian ρ -sampling of M . Let $m \in [0, 1[$ be a mass parameter such that $k = m|P|$, then the empirical measure μ on P is an (ϵ, r) -sampling of M for:*

$$\epsilon \geq \rho \sqrt{\frac{2 + k^2}{3}}, \quad \text{and } r = \infty$$

Proof: We consider the worst case scenario. Let M' be the segment centred at 0 and of length $2\rho n$ on the real line embedded in \mathbb{R}^d . We construct the sparsest possible ρ -sampling of M' . This is obtained by putting points at regular intervals 2ρ along the real line. We consider the empirical measure μ' on P' .

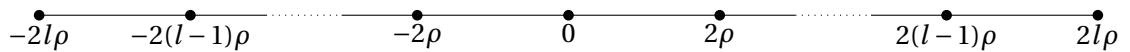


Figure 6.5 – The set P' sampling the segment $[-2l\rho, 2l\rho]$ when $n = 2l + 1$.

First we prove that M' is a worst case scenario for the parameter ϵ . It means that for any $x \in M$, we should be able to find a $x' \in M'$ such that $d_{\mu', m}(x') \geq d_{\mu, m}(x)$. Let us consider a point $x \in M$.

Then the nearest neighbour of x in P is $p_1(x)$. By hypothesis on P , $d_M(x, p_1(x)) \leq \rho$. Consider now the point x' of coordinate $d_M(x, p_1(x))$. Then the nearest neighbour of x' is 0 and it is at distance exactly $d_M(x, p_1(x))$.

We want to prove that $|P' \cap B(x', r)| \leq |P \cap B(x, r)|$ for any r . We built x' such that $|P' \cap B(x', r)| = 0$ if $r < d_M(x, p_1(x))$. Let us now assume that $|P' \cap B(x', r)| = i$. This implies that $r > (i-1)\rho + d_M(x, p_1(x))$ if i is odd and $r > i\rho - d_M(x, p_1(x))$ if i is even.

We only prove the result for i odd. Let us consider the Euclidean ball of radius $r = (i-1)\rho + d_M(x, p_1(x))$ and centre x . This Euclidean ball contains the Riemannian ball \mathcal{B} with the same radius and centre. Either this ball contains all M and in this case $|P \cap B(x, r)| = n = |P'|$, or the ball contains two points a and b such that the minimizing geodesic between a and b is enclosed in the ball and has length $2r$, because M is connected. Given that $r > (i-1)\rho + d_M(x, p_1(x))$, we can build a set of $i-1$ disjoint Riemannian balls of radius ρ included in \mathcal{B} such that none contains the point $p_1(x)$. For example, start by putting balls tangent in $p_1(x)$. By the ρ -sampling hypothesis each one of these balls must contain a point of P and hence \mathcal{B} contains at least i points which gives us the result.

A similar reasoning holds for the case i even.

Now, we compute the maximum of $d_{\mu', m}$ over the interval $[0, \rho]$. Let us fix $\alpha \in [0, 1]$.

$$\begin{aligned} d_{\mu', m}(\alpha\rho)^2 &= \frac{1}{k} \left(\sum_{i=0}^{\lfloor \frac{k-1}{2} \rfloor} (\alpha + 2i)^2 \rho^2 + \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor - 1} (2i + 2 - \alpha)^2 \rho^2 \right) \\ &= \frac{\rho^2}{k} \left(\sum_{i=0}^{\lfloor \frac{k-1}{2} \rfloor} (\alpha + 2i)^2 + \sum_{i=1}^{\lfloor \frac{k}{2} \rfloor} (2i - \alpha)^2 \right) \end{aligned}$$

Case k even : Fix $k = 2l$. Then :

$$\begin{aligned} d_{\mu', m}(\alpha\rho)^2 &= \frac{\rho^2}{2l} \left(\alpha^2 + \sum_{i=1}^{l-1} (2\alpha^2 + 8i^2) + (2l - \alpha)^2 \right) \\ &= \frac{\rho^2}{2l} \left(2l\alpha^2 - 4l\alpha + 4l^2 + 8 \sum_{i=1}^{l-1} i^2 \right) \\ &= \frac{\rho^2}{2l} \left(2l\alpha^2 - 4l\alpha + 4l^2 + \frac{8}{6}(l-1)l(2l-1) \right) \end{aligned}$$

The maximum of this expression is reached for $\alpha = 0$. Hence,

$$d_{\mu', m}(\alpha\rho)^2 \leq \frac{2\rho^2}{3}(2l^2 + 1)$$

Case k odd : Fix $k = 2l + 1$. Then:

$$d_{\mu', m}(\alpha\rho)^2 = \frac{\rho^2}{2l+1} \left((2l+1)\alpha^2 + 8 \sum_{i=1}^l i^2 \right)$$

The maximum is thus obtained for $\alpha = 1$. Hence,

$$\begin{aligned} d_{\mu',m}(\alpha\rho)^2 &\leq \frac{\rho^2}{2l+1} \left((2l+1) + 8 \frac{l(l+1)(2l+1)}{6} \right) \\ &\leq \frac{\rho^2}{3} (3 + 4l(l+1)) \\ &\leq \frac{\rho^2}{3} ((2l+1)^2 + 2) \end{aligned}$$

When we combine both results, we obtain the general relation for $k \geq 1$:

$$d_{\mu,m}(x) \leq d_{\mu',m}(\alpha\rho)^2 \leq \frac{\rho^2}{3} (2 + k^2)$$

Thus we conclude :

$$\epsilon \leq \rho \sqrt{\frac{2 + k^2}{3}}$$

Remark that our sampling condition is without noise. Thus all points of P are on M and hence $d_{\mu,m}(x) \geq d(x, M)$ for any x in the ambient space, which implies $r = \infty$. ■

Remark that some of the hypotheses are clearly not necessary. The assumption that M is connected is not needed. In fact, we only need an upper bound of $d_{\mu,m}$ on M . This can be obtained as soon as each of the connected components of M is sampled with at least k points. Moreover, the value of ϵ can be improved depending on the dimension of M . These results do not interest us here as we only want to show that the notion of (ϵ, r) -sampling is relevant and encompasses the previous noise conditions.

If we assume that points can move from their original position within a bounded range σ , we still obtain an Euclidean ρ' -sampling of M , where $\rho' = \rho + \sigma$. The same construction gives an (ϵ, r) -sampling, where $\epsilon = \rho' \sqrt{\frac{2+k^2}{3}}$ and $r = \infty$. In fact, remark that $d(x, M) \leq \sigma + d(x, P) \leq \epsilon + d_{\mu,m}(x)$ for any $x \in \mathbb{R}^d$.

6.4 Consequences on persistence diagrams

The (ϵ, r) -sampling hypothesis yields partial guarantees on the persistence diagram. ϵ is a parameter giving the precision of the approximation we obtain, while r indicates the range at which we can see the data. Intuitively, this means that as we look at the data at a larger scale, the noise, for example coming from a scatter noise, will completely mask the useful information.

Consider a compact Riemannian submanifold $M \subset \mathbb{R}^d$ and μ an (ϵ, r) -sample of M . If $r < \infty$, the interleaving needed to apply Theorem 2.23 will not be obtained for all values of the filtration parameter. However, we can get some guarantees on parts of the persistence diagrams.

Definition 6.9 *Given a filtration $\{F_\alpha\}$, where F_α is a subset of \mathbb{R}^d for all α , and given $\delta \in \mathbb{R}$, we define the δ -collapsed filtration $\{\tilde{F}_\alpha\}$ such that $\tilde{F}_\alpha = F_\alpha$ for all $\alpha < \delta$ and $\tilde{F}_\alpha = \mathbb{R}^d$ otherwise.*

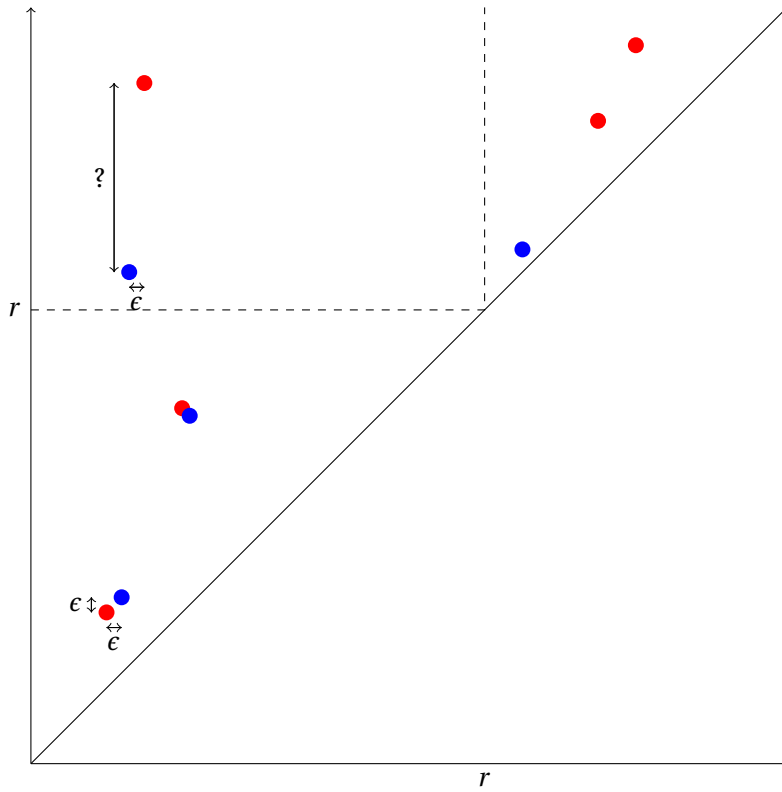


Figure 6.6 – Partition of the persistence diagrams according to the theoretical guarantees

Theorem 6.10 *Let μ be an (ϵ, r) -sample of a compact Riemannian submanifold $M \subset \mathbb{R}^d$. The $(r + \epsilon)$ -collapsed filtrations $\tilde{\mathcal{F}} = \{\tilde{F}_\alpha\}$ of $\{d_M^{-1}([0, \alpha])\}$ and $\tilde{\mathcal{G}} = \{\tilde{G}_\alpha\}$ of $\{d_{\mu, m}^{-1}([0, \alpha])\}$ are ϵ -interleaved and therefore:*

$$d_b(\text{Dgm}(\tilde{\mathcal{F}}), \text{Dgm}(\tilde{\mathcal{G}})) \leq \epsilon.$$

Proof: Consider $\alpha < r$. Then Lemma 6.2 implies that $\tilde{F}_\alpha \subset \tilde{G}_{\alpha+\epsilon}$ and $\tilde{G}_\alpha \subset \tilde{F}_{\alpha+\epsilon}$. Moreover, if $\alpha \geq r$ then $\tilde{F}_{\alpha+\epsilon} = \tilde{G}_{\alpha+\epsilon} = \mathbb{R}^d$. Therefore $\tilde{F}_\alpha \subset \tilde{G}_{\alpha+\epsilon}$ and $\tilde{G}_\alpha \subset \tilde{F}_{\alpha+\epsilon}$. Hence $\tilde{\mathcal{F}}$ and $\tilde{\mathcal{G}}$ are ϵ -interleaved and Corollary 2.27 applies. ■

Following the first steps of the proof of [26, Theorem 5.1], the persistence diagram $\text{Dgm}(\tilde{\mathcal{F}})$ of the r -collapsed filtration is identical to the persistence diagram of the original filtration $\text{Dgm}(\mathcal{F})$ in the quadrant $]-\infty, r[\times]-\infty, r[$. All points of $\text{Dgm}(\mathcal{F})$ located in $]-\infty, r[\times [r, \infty[$ are projected onto the line $\mathbb{R} \times r$. Finally, $\text{Dgm}(\mathcal{F})$ does not contain any element in $[r, \infty[\times [r, \infty[$. Theorem 6.10 therefore gives guarantees on parts of the persistence diagrams $\text{Dgm}(d_{\mu, m})$ and $\text{Dgm}(d_M)$. There exists a partial matching π of $\text{Dgm}(d_{\mu, m})$ with $\text{Dgm}(d_M)$ verifying the following conditions, where Δ is the diagonal.

1. $\forall x \in [0, r] \times [0, r] \cap \text{Dgm}(d_{\mu, m}), d_\infty(x, \Delta) \leq \epsilon \implies x$ is matched with $\pi(x)$ such that $\|\pi(x) - x\|_\infty \leq \epsilon$.
2. $\forall x = (a, b) \in [0, r] \times [r, \infty] \cap \text{Dgm}(d_{\mu, m}), d_\infty(x, \Delta) \leq \epsilon \implies x$ is matched with $\pi(x) = (a', b')$

such that $|a' - a| \leq \epsilon$.

The conditions are also valid if we reverse the roles of $\text{Dgm}(d_{\mu,m})$ and $\text{Dgm}(d_M)$. Intuitively, it means that the two diagrams are at distance at most ϵ in the bottom left part $[0, r] \times [0, r]$, that the difference between them in the upper quadrant $[0, r] \times]r, \infty]$ is only bounded horizontally and nothing is known for the remaining upper right part.

7 Scalar field analysis

In [32], Chazal et al. presented an algorithm to analyse the topology of a scalar field using persistent homology which can handle bounded Hausdorff noise both in geometry and in observed function values. We build upon the same framework in order to handle unbounded noise using the distance to a measure.

First, we introduce necessary preliminaries as well as some of the results from [32]. Then, we show how to handle unbounded functional and geometric noise and present some experimental illustration.

7.1 Scalar field analysis with bounded noise

Given a Riemannian manifold M , the scalar field topology of $f : M \rightarrow \mathbb{R}$ is studied via the topological structure of the sub-level sets filtration of f , defined as $F_\alpha = f^{-1}([-\infty, \alpha])$ for any $\alpha \in \mathbb{R}$. The collection of sub-level sets form a filtration $\mathcal{F} = \{F_\alpha\}_{\alpha \in \mathbb{R}}$ connected by natural inclusions $F_\alpha \subseteq F_\beta$ for any $\alpha \leq \beta$. Our goal is to approximate the persistence diagram $\text{Dgm}(\mathcal{F})$ from the observed scalar field $\tilde{f} : P \rightarrow \mathbb{R}$, where P is a set of measure points. We first describe the results of [32] for approximating $\text{Dgm}(\mathcal{F})$ when P is a geodesic ε -sampling of M .

To simulate the sub-level sets filtration $\{F_\alpha\}$ of f , we introduce $P_\alpha = \tilde{f}^{-1}([-\infty, \alpha]) \subset P$ for any $\alpha \in \mathbb{R}$. Intuitively, the points in P_α sample the sub-level set F_α . To estimate the topology of F_α from these discrete samples P_α , we consider the δ -offset P_α^δ of the point set P_α i.e. we grow balls of radius δ around the points of P_α . It gives us a union of balls that serves as a proxy for $f^{-1}([-\infty, \alpha])$ and whose nerve is the Čech complex, $C_\delta(P_\alpha)$. For computation purpose, we use the Vietoris-Rips complex $R_\delta(P_\alpha)$, which is related to $C_\delta(P_\alpha)$ as shown in Proposition 5.3: $\forall \delta > 0, C_\delta(P_\alpha) \subset R_\delta(P_\alpha) \subset C_{2\delta}(P_\alpha)$.

Even though no Vietoris-Rips complex might capture the topology of the manifold M , it was shown in [35] that a structure of nested complexes can recover it from the filtration $\{P_\alpha\}$, using the inclusions $R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)$. Specifically, for a fixed $\delta > 0$ and $\alpha \leq \beta$, consider the following commutative diagram induced by inclusions,

$$\begin{array}{ccc} H_*(R_{2\delta}(P_\alpha)) & \longrightarrow & H_*(R_{2\delta}(P_\beta)) \\ \uparrow & & \uparrow \\ H_*(R_\delta(P_\alpha)) & \longrightarrow & H_*(R_\delta(P_\beta)) \end{array}$$

Defining Φ_α as the image of $H_*(R_\delta(P_\alpha)) \rightarrow H_*(R_{2\delta}(P_\alpha))$, the diagram induces a map $\phi_\alpha^\beta : \Phi_\alpha \rightarrow$

Φ_β . As the diagram commutes for all $\alpha \leq \beta$, $\{\Phi_\alpha, \phi_\alpha^\beta\}$ defines a persistence module. We call it the persistent homology module of the filtration of nested pairs $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}_{\alpha \in \mathbb{R}}$.

Using this construction, one of the main results of [32] is:

Theorem 7.1 (Theorems 2 and 6 of [32]) *Let M be a compact Riemannian manifold and let $f : M \rightarrow \mathbb{R}$ be a c -Lipschitz function. Let P be a geodesic ε -sampling of M . If $\varepsilon < \frac{1}{4}\rho(M)$, then for any $\delta \in [2\varepsilon, \frac{1}{2}\rho(M)]$, the persistent homology modules of f and of the filtration of nested pairs $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$ are $2c\delta$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most $2c\delta$.*

Furthermore, the k -dimensional persistence diagram for the filtrations of nested pairs $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$ can be computed in $O(|P|kN + N \log N + N^3)$ time, where N is the number of simplices of $\{R_{2\delta}(P_\infty)\}$, and $|P|$ denotes the cardinality of the sample set P .

It has been observed that in practice, the persistence algorithm often has a running time linear in the number of simplices, which reduces the above complexity to $O(|P| + N \log N)$ in a practical setting.

We say that \tilde{f} has a precision of ξ over P if $|\tilde{f}(p) - f(p)| \leq \xi$ for any $p \in P$. This is what we call a Hausdorff type functional noise and we have:

Theorem 7.2 (Theorem 3 of [32]) *Let M be a compact Riemannian manifold and let $f : M \rightarrow \mathbb{R}$ be a c -Lipschitz function. Let P be a geodesic ε -sampling of M such that the values of f on P are known with precision ξ . If $\varepsilon < \frac{1}{4}\rho(M)$, then for any $\delta \in [2\varepsilon, \frac{1}{2}\rho(M)]$, the persistent homology modules of f and of the filtration of nested pairs $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$ are $(2c\delta + \xi)$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most $2c\delta + \xi$.*

Geometric noise was considered in the form of bounded noise in the estimate of the geodesic distances between points in P . It translated into a relation between the measured pairwise distances and the real ones to obtain:

Theorem 7.3 (Theorem 4 of [32]) ¹*Let M , f be defined as previously and P be an ε -sample of M in its Riemannian metric. Assume that, for a parameter $\delta > 0$, the Rips complexes $R_\delta(\cdot)$ are defined with respect to a metric $\tilde{d}(\cdot, \cdot)$ which satisfies $\forall x, y \in P, \frac{d_M(x, y)}{\lambda} \leq \tilde{d}(x, y) \leq v + \mu \frac{d_M(x, y)}{\lambda}$, where $\lambda \geq 1$ is a scaling factor, $\mu \geq 1$ is a relative error and $v \geq 0$ an additive error. Then, for any $\delta \geq v + 2\mu\frac{\varepsilon}{\lambda}$ and any $\delta' \in [v + 2\mu\delta, \frac{1}{\lambda}\rho(M)]$, the persistent homology modules of f and of the filtration of nested pairs $\{R_\delta(P_\alpha) \hookrightarrow R_{\delta'}(P_\alpha)\}$ are $c\lambda\delta'$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most $c\lambda\delta'$.*

7.2 Unbounded functional noise

In this section, we focus on the case where we only have noise in the observed function \tilde{f} . Suppose that we have a scalar function f defined on a manifold M embedded in a metric space \mathbb{X} . We are given a geodesic ε -sample $P \subset M$, and a noisy observed function $\tilde{f} : P \rightarrow \mathbb{R}$. Our goal is to approximate the persistence diagram $\text{Dgm}(\mathcal{F})$ of the sub-level sets filtration

¹This is a specific case of Theorem 4 of [32] and was stated in this form in the conference version of the paper.

$\mathcal{F} = \{F_\alpha = f^{-1}((-\infty, \alpha])\}_\alpha$ from \tilde{f} . We assume that f is c -Lipschitz with respect to the intrinsic metric of the manifold M . Note that this does not imply a Lipschitz condition on f for the extrinsic metric.

7.2.1 Functional noise model

Previous work on functional noise usually focuses on Hausdorff-type bounded noise [32] or statistical noise with zero-mean [74]. However, we observe that there are many practical scenarios where the observed function \tilde{f} may contain these previously considered types of noise mixed with *aberrant function values*. Hence, we propose below a more general noise model that allows such a mixture.

Motivating examples. First, we provide some motivating examples for the need of handling *aberrant* function values in \tilde{f} , i.e. where $\tilde{f}(p)$ at some sample p can be totally unrelated to the true value $f(p)$. Consider a sensor network, where each node returns some measures. Such measurements can be imprecise, and in addition to that, a sensor may experience failure and return a completely wrong measure that has no relation with the true value of f . Similarly, an image could be corrupted with impulse noise, where there are random pixels with aberrant function values, such as random white or black dots (see Figure 7.4 for an illustration).

More interestingly, outliers in function values can naturally appear as a result of geometric noise present in the discrete samples. For example, imagine that we have a process that can measure the function value $f : M \rightarrow \mathbb{R}$ with *no error*. However, the geometric location \tilde{p} of a point $p \in M$ can be wrong. In particular, \tilde{p} can be close to other parts of the manifold, thereby although \tilde{p} has the correct function value $f(p)$, it becomes a functional outlier among its neighbours, due to the wrong location of \tilde{p} . Figure 7.1 shows a bone-like structure, where the function $f(x)$ is given by the distance on the curve between x and the point located at the middle of the bottom horizontal segment. The two sides of the narrow neck have very different function values. Consider that the points are sampled uniformly on M and their position is then convolved with a Gaussian noise in the ambient space. Then points from one side of this neck can be sent closer to the other side, causing aberrant values in the observed function.

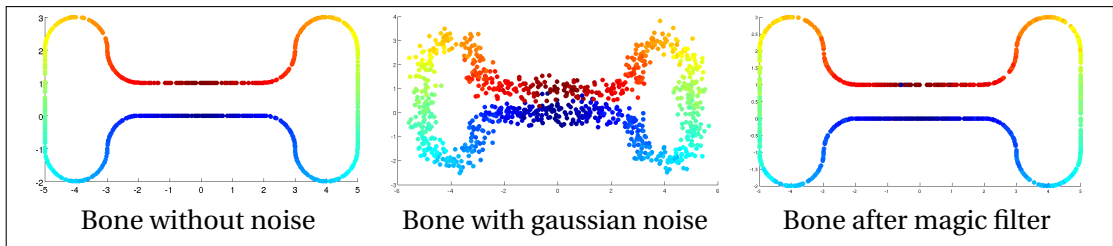


Figure 7.1 – Bone example after applying Gaussian perturbation and magic filter

In fact, if we assume that we have a “magic filter” that can project each sample back onto the underlying manifold M , thus obtaining a new set of samples where all points are on the manifold, we obtain a sampling that can be seen as having no geometric noise, but a functional noise caused by the original geometric noise. Note that such a magic filter is the goal of many

geometric denoising methods. This implies that a denoising algorithm, perfect in the sense of geometric noise, cannot remove or may even cause more aberrant functional noise. This motivates the need for handling functional outliers in addition to traditional functional noise as well as processing noise that combines geometric and functional noise together and is not necessarily centred. Figure 7.1 shows a bone-like curve and a function defined as the curvilinear abscissa. The Gaussian noise applied to the example creates outliers even after applying a projection onto the original object.

Another case where our approach is useful concerns missing data. Assuming that some of the functional values are missing, we can replace them by anything and act as if they were outliers. Without modifying the algorithm, we obtain a way to handle the local loss of information. This will be further discussed in Chapter 8.

Functional noise model. To allow both aberrant and more traditional functional noise, we introduce a new noise model. Given a set of points P and an integer k , we write $\text{NN}_P^k(p)$ the set of the k -nearest neighbours of p in P in the *extrinsic metric*.

Definition 7.4 Let M be a Riemannian manifold in a metric space \mathbb{X} and $f : M \rightarrow \mathbb{R}$ a scalar field on M . Given a geodesic ϵ -sample P of M , a discrete scalar field $\tilde{f} : P \rightarrow \mathbb{R}$ and two integers k and k' such that $k \geq k' > \frac{k}{2}$, we say that \tilde{f} is a (k, k', Δ) -functional-sample of f if

$$\forall p \in P, \left| \left\{ q \in \text{NN}_P^k(p) \mid |\tilde{f}(q) - f(p)| \leq \Delta \right\} \right| \geq k' \quad (7.1)$$

Intuitively, this noise model allows up to $k - k'$ samples around a point p to be outliers, i.e. points whose functional value deviates from $f(p)$ by at least Δ . We now consider several common functional noise models used in the statistical learning community and look at what they correspond to in our setting.

Bounded noise model. The standard “bounded noise” model assumes that all observed function values are within some distance δ away from the true function values, that is, $|\tilde{f}(p) - f(p)| \leq \delta$ for all $p \in P$. Hence this bounded noise model simply corresponds to a $(1, 1, \delta)$ -functional-sample.

Gaussian noise model. Under the popular Gaussian noise model, for any $x \in M$, its observed function value $\tilde{f}(x)$ is drawn from a normal distribution $\mathcal{N}(f(x), \sigma)$, that is a probability measure with density $g(y) = \frac{1}{\sigma\sqrt{\pi}} e^{-\frac{(y-f(x))^2}{\sigma^2}}$. We say that a point $q \in P$ is a -accurate if $|\tilde{f}(q) - f(q)| \leq a$. For the Gaussian noise model, we will first bound the random variable $\mu(k, k')$ defined as the smallest value such that at least k' out of the k nearest neighbours of p in $\text{NN}_P^k(p)$ are $\mu(k, k')$ -accurate.

Lemma 7.5 With probability at least $1 - e^{-\frac{k-k'}{6}}$, $\mu(k, k') \leq \sigma \sqrt{\ln \frac{2k}{k-k'}}$.

Proof: First note that $\frac{1}{\sqrt{2\pi}} \int_b^{+\infty} e^{-\frac{t^2}{2}} dt \leq \frac{1}{\sqrt{2\pi}} e^{-\frac{b^2}{2}}$ if $b \geq 1$. Introducing $I(a) = \frac{1}{\sigma\sqrt{\pi}} \int_{-a}^a e^{-\frac{x^2}{\sigma^2}} dx$, we get, by a change of variable, for $a \geq \sigma$:

$$1 - e^{-\left(\frac{a}{\sigma}\right)^2} \leq I(a) \quad (7.2)$$

Now set $\delta = \frac{k-k'}{k} < \frac{1}{2}$ and $s = \sigma \sqrt{\ln \frac{2k}{k-k'}} \geq \sigma$. Let p_1, \dots, p_k denote the k nearest neighbours of some point, say p_1 . For each p_i , let $Z_i = 1$ if p_i is **not** s -accurate, and $Z_i = 0$ otherwise. Hence $Z = \sum_{i=1}^k Z_i$ denotes the total number of points from these k nearest neighbours that are not s -accurate. By Equation (7.2), we know that

$$\text{Prob}[Z_i = 1] = 1 - I(s) \leq e^{-\left(\frac{s}{\sigma}\right)^2}.$$

It then follows that the expected value of Z satisfies:

$$\mathbb{E}(Z) \leq k e^{-\left(\frac{s}{\sigma}\right)^2} = \frac{\delta k}{2}.$$

Now set $\rho = \frac{\delta k}{2\mathbb{E}(Z)}$. Since $\mathbb{E}(Z) \leq \frac{\delta k}{2}$, it follows that $(1 + \rho)\mathbb{E}(Z) \leq \delta k$. Using Chernoff's bound [5], we obtain

$$\begin{aligned} \text{Prob}[Z \geq k - k'] &= \text{Prob}[Z \geq \delta k] \leq \text{Prob}[Z \geq (1 + \rho)\mathbb{E}(Z)] \\ &\leq e^{-\frac{\rho^2 \mathbb{E}(Z)}{2 + \rho}} = e^{-\frac{\delta^2 k^2}{4\mathbb{E}(Z)} \frac{1}{2 + \frac{\delta k}{\mathbb{E}(Z)}}} \leq e^{-\frac{\delta^2 k^2}{6\delta k}} = e^{-\frac{k-k'}{6}}. \end{aligned}$$

We remark that, as we increase the value of s , the probability $\text{Prob}[Z \geq k - k']$ of having at least $k - k'$ points not s -accurate decreases. ■

Next, we convert the value $\mu(k, k')$ to the value Δ as in Equation (7.1).

Proposition 7.6 *Let $f : M \rightarrow \mathbb{R}$ be a c -Lipschitz scalar field on a Riemannian manifold M . Given a sampling P of M and a discrete scalar field $\tilde{f} : P \rightarrow \mathbb{R}$ respecting the Gaussian noise model, then with probability at least $1 - ne^{-\frac{k-k'}{\sigma}}$, where $n = |P|$, \tilde{f} is a (k, k', Δ) -functional-sample for $\Delta = \sigma \sqrt{\ln \frac{2k}{k-k'}} + c\lambda$, where λ is the geodesic diameter of $\text{NN}_P^k(p)$.*

Proof: In particular, being a (k, k', Δ) -functional-sample means that for any $p \in P$, there are at least k' samples q from $\text{NN}_P^k(p)$ such that $|\tilde{f}(q) - f(p)| \leq \Delta$. Now assume that the furthest geodesic distance from any point in $\text{NN}_P^k(p)$ to p is λ . Then since f is a c -Lipschitz function, we have $\max_{q \in \text{NN}_P^k(p)} |f(q) - f(p)| \leq c\lambda$.

We note that Lemma 7.5 is valid for any point p of P . Using the union bound, the relation holds for all points in P with probability at least $1 - ne^{-\frac{k-k'}{6}}$. Note that if $k - k' \geq 12 \ln n$, then this probability is at least $1 - \frac{1}{n}$, that is, the relation holds with high probability. Thus, with probability at least $1 - ne^{-\frac{k-k'}{6}}$, the input function $\tilde{f} : P \rightarrow \mathbb{R}$ under Gaussian noise model is a (k, k', Δ) -functional-sample with $\Delta = \sigma \sqrt{\ln \frac{2k}{k-k'}} + c\lambda$. ■

7.2.2 Functional denoising

Given a scalar field $\tilde{f} : P \rightarrow \mathbb{R}$ which is a (k, k', Δ) -functional-sample of $f : M \rightarrow \mathbb{R}$, we now aim at computing a denoised function $\hat{f} : P \rightarrow \mathbb{R}$ from the observed function \tilde{f} , and we will later use \hat{f} to infer the topology of the sub-level sets of $f : M \rightarrow \mathbb{R}$. Below we describe two ways to denoise \tilde{f} . One is well-known while the other one is new. As we will see later, these two treatments lead to similar theoretical guarantees in terms of topology inference. However, they have different characteristics in practice, which we will discuss in the experimental illustration of Section 7.2.3.

k -median. In the k -median treatment, we simply perform the following: given any point $p \in P$, we set $\hat{f}(p)$ to be the median value of the set of \tilde{f} values for the k -nearest neighbours $\text{NN}_p^k(p) \subseteq P$ of p . We call \hat{f} the k -median denoising of \tilde{f} . The following observation is simple:

Lemma 7.7 *If $\tilde{f} : P \rightarrow \mathbb{R}$ is a (k, k', Δ) -functional-sample of $f : M \rightarrow \mathbb{R}$ with $k' \geq k/2$, then we have $|\hat{f}(p) - f(p)| \leq \Delta$ for any $p \in P$, where \hat{f} is the k -median denoising of \tilde{f} .*

Proof: Given $p \in P$, there exist at least k' points x_i among its k -nearest neighbours such that $f(x_i) \geq f(p) - \Delta$. Similarly, there exist at least k' points such that $f(y_i) \leq f(p) + \Delta$. Since $k' > \frac{k}{2}$, $f(p) - \Delta \leq \hat{f}(p) \leq f(p) + \Delta$. ■

Discrepancy. In the k -median treatment, we choose a single value from the k -nearest neighbours of a sample point p and set it to be the denoised value $\hat{f}(p)$. This value, while within distance Δ from the true value $f(p)$ when $k' \geq k/2$, tends to have greater variability among neighbouring sample points. Intuitively, taking the average, such as k -means, makes the function $\hat{f}(p)$ smoother, but it is sensitive to outliers.

The distance to a measure function we used in previous chapters provided an averaging that was less sensible to outliers. Considering the set N of the k nearest neighbours of a point, we can look at their functional values and consider the distance to the empirical measure on $\tilde{f}(N)$. We look for the minimum of this function and use it as our estimate $\hat{f}(p)$.

Given a set $Y = \{x_1, \dots, x_{k'}\}$ of k' sample points from P , we define its discrepancy with respect to \tilde{f} as:

$$\phi(Y) = \frac{1}{k'} \sum_{i=1}^{k'} (\tilde{f}(x_i) - \mu(Y))^2, \quad \text{where } \mu(Y) = \frac{1}{k'} \sum_{i=1}^{k'} \tilde{f}(x_i).$$

$\mu(Y)$ and $\phi(Y)$ are respectively the mean and the variance of the observed function values for points from Y . Intuitively, $\phi(Y)$ measures how tight the function values $(\tilde{f}(x_i))$ are clustered. Now, given a point $p \in P$, we define

$$\hat{Y}_p = \operatorname{argmin}_{Y \subseteq \text{NN}_p^k(p), |Y|=k'} \phi(Y), \quad \text{and } \hat{z}_p = \mu(\hat{Y}_p).$$

That is, \hat{Y}_p is the subset of k' points from the k -nearest neighbours of p that has the smallest discrepancy and \hat{z}_p is its mass centre. It turns out that \hat{Y}_p and \hat{z}_p can be computed by the following sliding-window procedure:

 SLIDING WINDOW PROCEDURE

1. Sort $\text{NN}_p^k(p) = \{x_1, \dots, x_k\}$ according to $\tilde{f}(x_i)$.
 2. For every k' consecutive points $Y_i = \{x_i, \dots, x_{i+k'-1}\}$ with $i \in [1, k - k' + 1]$, compute its discrepancy $\phi(Y_i)$.
 3. Set $\hat{Y}_p = \operatorname{argmin}_{Y_i, i \in [1, k - k']} \phi(Y_i)$, and return $\mu(\hat{Y}_p)$ as \hat{z}_p .
-

Remark that \hat{z}_p is exactly the minimum over \mathbb{R} of the distance to the empirical measure on $f(\text{NN}_p^k(p))$ for the mass $\frac{k'}{k}$. In the *discrepancy-based denoising*² approach, we simply set $\hat{f}(p) := \hat{z}_p$ as computed above. The correctness of \hat{f} to approximate f is given by the following Lemma.

Lemma 7.8 *If $\tilde{f} : P \rightarrow \mathbb{R}$ is a (k, k', Δ) -functional-sample of $f : M \rightarrow \mathbb{R}$ with $k' \geq \frac{k}{2}$, then we have $|\hat{f}(p) - f(p)| \leq \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right) \Delta$ for any $p \in P$, where \hat{f} is the discrepancy-based denoising of \tilde{f} . In particular, if $k' \geq \frac{2}{3}k$, then $|\hat{f}(p) - f(p)| \leq 3\Delta$ for any $p \in P$.*

Proof: Let $Y_\Delta = \{x \in \text{NN}_p^k(p) : |\tilde{f}(x) - f(p)| \leq \Delta\}$ be the set of points in $\text{NN}_p^k(p)$ whose observed function values are at distance at most Δ away from $f(p)$. Since \tilde{f} is a (k, k', Δ) -functional-sample of f , $|Y_\Delta| \geq k'$. Let $Y'_\Delta \subset Y_\Delta$ be a subset with k' elements, $Y'_\Delta = \{x'_i\}_{i=1}^{k'}$. By definition of Y_Δ and Y'_Δ , $|\tilde{f}(x'_i) - \mu(Y'_\Delta)| \leq 2\Delta$ where $\mu(Y'_\Delta) = \frac{1}{k'} \sum_{i=1}^{k'} \tilde{f}(x'_i)$. This inequality gives an upper bound of the discrepancy $\phi(Y'_\Delta)$,

$$\begin{aligned} \phi(Y'_\Delta) &= \frac{1}{k'} \sum_{i=1}^{k'} (\tilde{f}(x'_i) - \mu(Y'_\Delta))^2 \\ &\leq \frac{1}{k'} \sum_{i=1}^{k'} (2\Delta)^2 \\ &= 4\Delta^2 \end{aligned}$$

Recall that $\hat{Y}_p = \operatorname{argmin}_{Y_i, i \in [1, k - k']} \phi(Y_i)$ and $\hat{z}_p = \mu(\hat{Y}_p)$. Denote $A_1 = \hat{Y}_p \cap Y_\Delta$ and $A_2 = \hat{Y}_p \setminus A_1$. Since \tilde{f} is a (k, k', Δ) -functional-sample of f , the size of A_2 is at most $k - k'$ and $|A_1| \geq 2k' - k$. If $|\hat{z}_p - f(p)| \leq \Delta$, nothing needs to be proved. Without loss of generality, we assume that $f(p) + \Delta \leq \hat{z}_p$. Denote $\delta = \hat{z}_p - (f(p) + \Delta)$. The discrepancy of $\phi(\hat{Y}_p)$ can be estimated as follows.

²It has no relation to the discrepancy method [36].

$$\begin{aligned}
 \phi(\hat{Y}_p) &= \frac{1}{k'} \left(\sum_{x \in A_1} (\tilde{f}(x) - \hat{z}_p)^2 + \sum_{x \in A_2} (\tilde{f}(x) - \hat{z}_p)^2 \right) \\
 &\geq \frac{1}{k'} \left(|A_1| \delta^2 + \sum_{x \in A_2} (\tilde{f}(x) - \hat{z}_p)^2 \right) \\
 &\geq \frac{1}{k'} \left(|A_1| \delta^2 + \frac{1}{|A_2|} \left(\sum_{x \in A_2} \tilde{f}(x) - |A_2| \hat{z}_p \right)^2 \right) \\
 &= \frac{1}{k'} \left(|A_1| \delta^2 + \frac{1}{|A_2|} \left(\sum_{x \in A_1} \tilde{f}(x) - |A_1| \hat{z}_p \right)^2 \right) \\
 &\geq \frac{1}{k'} \left(|A_1| \delta^2 + \frac{1}{|A_2|} (|A_1| \delta)^2 \right) \\
 &\geq \frac{1}{k'} \delta^2 \left(\frac{k' |A_1|}{|A_2|} \right) \\
 &\geq \frac{2k' - k}{k - k'} \delta^2
 \end{aligned}$$

where the third line uses the inequality $\sum_{i=1}^n a_i^2 \geq \frac{1}{n} (\sum_{i=1}^n a_i)^2$, and the fourth line uses the fact that $(|A_1| + |A_2|) \hat{z}_p = \sum_{x \in \hat{Y}_p} \tilde{f}(x)$. Since $\hat{Y}_p = \operatorname{argmin}_{Y_i, i \in [1, k-k']} \phi(Y_i)$, $\phi(\hat{Y}_p) \leq \phi(Y'_\Delta)$. Therefore,

$$\frac{2k' - k}{k - k'} \delta^2 \leq 4\Delta^2.$$

Hence $\delta \leq 2\sqrt{\frac{k-k'}{2k'-k}} \Delta$ and $|\hat{f}(p) - f(p)| \leq \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right) \Delta$, since $\hat{z}_p = \hat{f}(p)$. If $k' \geq \frac{2}{3}k$, then $1 + 2\sqrt{\frac{k-k'}{2k'-k}} \leq 1 + 2 = 3$, meaning that $|\hat{f}(p) - f(p)| \leq 3\Delta$. ■

Corollary 7.9 *Given a (k, k', Δ) -functional-sample of $f : M \rightarrow \mathbb{R}$ with $k' \geq k/2$, we can compute a new function $\hat{f} : P \rightarrow \mathbb{R}$ such that $|\hat{f}(p) - f(p)| \leq \xi \Delta$ for any $p \in P$, where $\xi = 1$ under k -median denoising, and $\xi = \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right)$ under the discrepancy-based denoising.*

After the k -median denoising or the discrepancy-based denoising, we obtain a new function \hat{f} whose value at each sample point has an error at most $\xi \Delta$ from the true function value. We can now apply the scalar field topology inference framework from [32], as introduced in Section 7.1, using \hat{f} as input. In particular, set $L_\alpha = \{p \in P \mid \hat{f}(p) \leq \alpha\}$, and let $R_\delta(X)$ denote the Rips complex over points in X with parameter δ . We approximate the persistence diagram induced by the sub-level sets filtration of $f : M \rightarrow \mathbb{R}$ from the filtration of nested pairs $\{R_\delta(L_\alpha) \hookrightarrow R_{2\delta}(L_\alpha)\}_\alpha$. It follows from Theorem 7.2 that:

Theorem 7.10 *Let M be a compact Riemannian manifold and let $f : M \rightarrow \mathbb{R}$ be a c -Lipschitz function. Let P be a geodesic ε -sampling of M , and $\hat{f} : P \rightarrow \mathbb{R}$ a (k, k', Δ) -functional-sample of f . Set $\xi = 1$ if P_α is obtained via k -median denoising, and $\xi = \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right)$ if P_α is obtained*

via discrepancy-based denoising. If $\varepsilon < \frac{1}{4}\varrho(M)$, then for any $\delta \in [2\varepsilon, \frac{1}{2}\varrho(M))$, the persistent homology modules of f and the filtration of nested pairs $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$ are $(2c\delta + \xi\Delta)$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most $2c\delta + \xi\Delta$.

Theoretical results are similar for k -median and discrepancy-based methods with a slight advantage for the k -median. However, interesting experimental results can be obtained when the Lipschitz condition on f is removed, for example with images, where the discrepancy based method shows more robustness to large amounts of noise. Moreover, the approach through the distance to a measure can be defined for other spaces of functional values than \mathbb{R} . In fact, it can be defined for any metric space \mathbb{X} , but the computation of the minimum is a challenge. The rest of the section presents an illustration of the behaviour of these methods.

7.2.3 Experimental Illustration for functional noise

In this section, we present results obtained by applying our methods to cases where there is only functional noise. Our goals are to demonstrate the denoising power of both the k -median and the discrepancy-based approaches and to illustrate the differences between the practical performances of the k -median and discrepancy-based denoising methods. We compare our denoising results with the popular k -NN algorithm, which simply sets the function at point p to be the mean of the observed function values of its k nearest neighbours. Note that, when $k' = k$, our discrepancy-based method is equivalent to the k -NN algorithm.

Going back to the bone example from section 7.2.1, we apply our algorithm to the 10-nearest neighbours and $k' = 8$. Using 100 sampling of the Bone with 1000 points each, we compute the average maximal error made by the various methods. The discrepancy-based method commits a maximal error of 10% on average, while the median-based method recovers the values with an error of 2% and the simple k -NN regression gives a maximal error of 16%, with most error concentrated around the neck region, see Figure 7.2. These results translate into the persistence diagrams that are more robust with the use of the discrepancy (blue squares) or the k -median (red diamond) instead of the k -NN regression (green circles), see Figure 7.3. Both methods retrieve the 1-dimensional topological feature. The k -NN regression keeps some prominent 0-dimensional feature through the diagram instead of having a unique component, result obtained by using the discrepancy or the median. The persistence diagram of the original bone is given in red and contains only one feature.

As indicated by theoretical results, the discrepancy-based method improves the classic k -NN regression but the median-based algorithm is still better. The discrepancy can however have a much better behaviour when the noise can have values in between two correct ones and when the Lipschitz condition is relaxed. This is the case in some other practical applications like image denoising.

We take the greyscale image Lena as the target scalar field f . In Figure 7.4, we use two ways to generate a noisy input scalar field \tilde{f} . The first type of noisy input is generated by adding uniform random noise, also called impulse noise, as follows: with probability p , each pixel will receive a uniformly distributed random value in range $[0, 255]$ as its function value; otherwise, it is unchanged. Results under random noises are in the second and third rows of Figure 7.4.

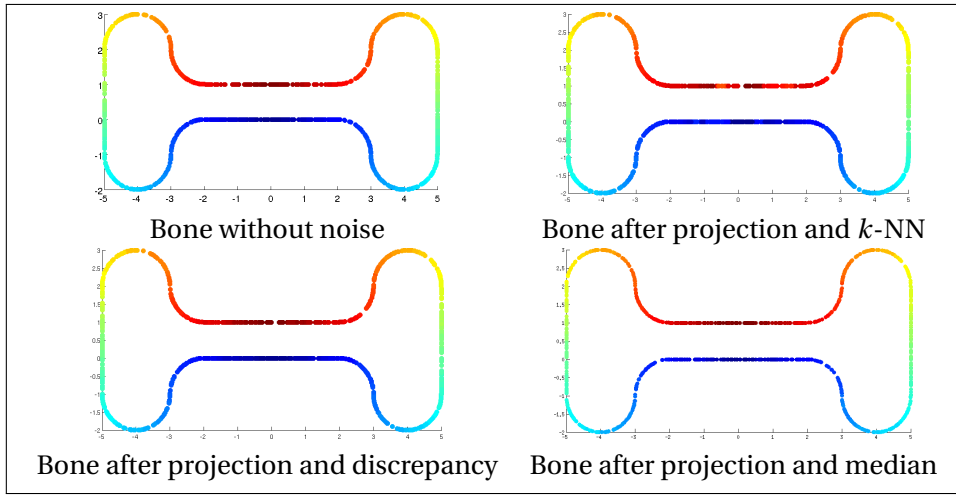


Figure 7.2 – Bone example after applying Gaussian perturbation, magical filter and a regression

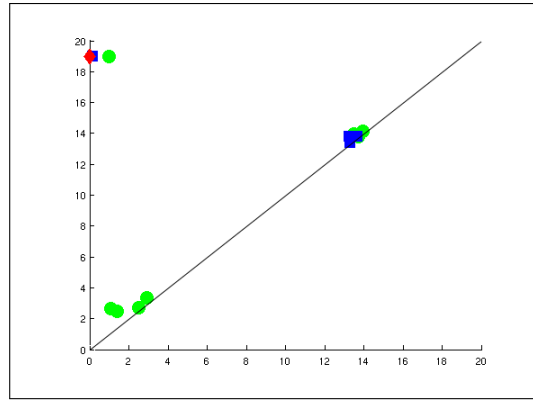


Figure 7.3 – Persistence diagrams in dimension 0 for the Bone example

We also consider what we call *outlier noise*, similar to the so-called salt and pepper noise: with probability p , each pixel will be an outlier meaning that its function value is a fixed constant, which is set to be 200 in our experiments. This outlier noise is to simulate the aberrant function values caused by a broken sensor for example. The denoising results under the outlier-noise are shown in the last row of Figure 7.4.

First, we note that the k -NN approach tends to smooth out function values. In addition to the blurring artifact, its denoising capability is limited when the amount of noise is high i.e. where imprecise values become dominant. As expected, both k -median and discrepancy based methods outperform the k -NN approach. Indeed, they demonstrate robust recovery of the input image even with 50% amount of random noise are added.

While both k -median and discrepancy based methods are more resilient against noise, there are interesting differences between their practical performances. From a theoretical point of

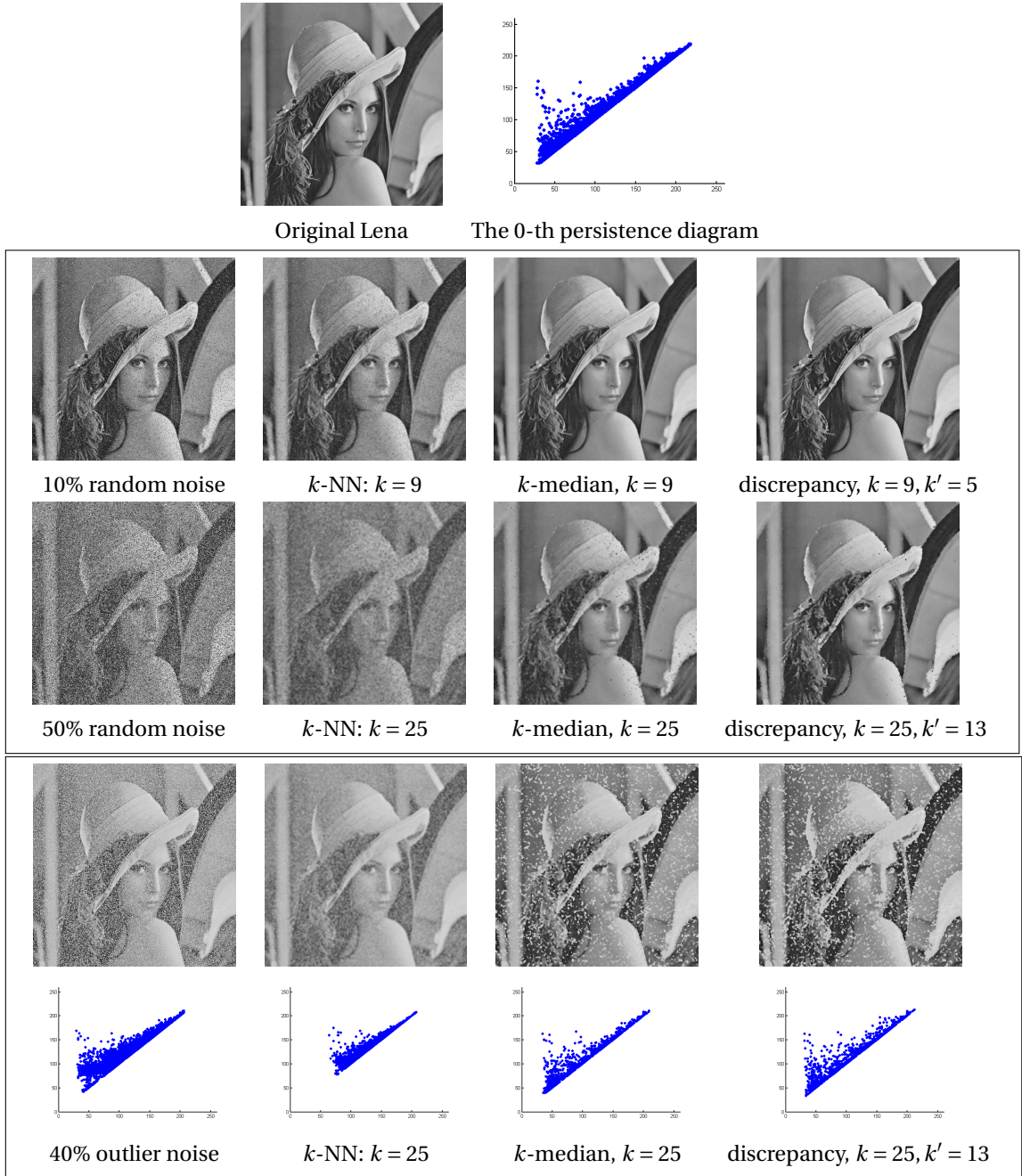


Figure 7.4 – The denoised images after k -NN, k -median, and discrepancy denoising approaches. The first row shows the original image and its 0-th persistence diagram. Second and third rows are under random noise of input, while fourth row are under outlier-noise as described in the text. The fifth row provides the 0-th persistence diagrams on images in the fourth row, which are computed by the scalar field analysis algorithm from [32].

view, when the input scalar field is indeed a (k, k', Δ) -functional-sample, the k -median method gives a slightly better error bound than the discrepancy based method, see Theorem 7.10. However, when the (k, k', Δ) -functional-sample condition is not satisfied, the median value can be quite arbitrary. By taking the average of a subset of points, the discrepancy method is more robust against large amount of noise. This difference is evident in the third and last row of Figure 7.4. In this case, the addition of noise creates a bias in the functional values. When the proportion of noise is sufficiently high and the values of noisy points are in between two groups of relevant values, for example, when you have a dark and clear parts in an image and the noise is grey, the median method returns a value in the grey, which is uncorrelated to relevant information, while the discrepancy-based method can remain closer to the original information. The main reason for this behaviour comes from the absence of the Lipschitz property in the image application. More precisely, the Lipschitz constant is too big to provide meaningful guarantees on the quality of the regression. The same problem also occurs in presence of random impulse noise. With a high proportion of noise, the median is very often close to the median of all possible values, while the discrepancy-based method returns a value close to denser region around relevant information.

Moreover, the application to persistent homology, which was our primary, goal is much cleaner after the discrepancy-based method. The structure of the beginning of the diagrams is almost perfectly retrieved by both the median and discrepancy-based methods. However, the median induces a shrinking phenomenon to the diagram. This means that the width of the diagram is reduced and so are the lifespans of topological features, making it more difficult to distinguish between noise and relevant information. Remark that the classic k -NN approach shrinks the diagram even more, to the point where it is no longer possible to distinguish the information from the noise.

The standard indicator to measure the quality of a denoising is the *Peak Signal over Noise Ratio* (PSNR). Given a grey scale input image I and an output image O with the grey scale between 0 and 255, it is defined by

$$\text{PSNR}(I, O) = 10 \log_{10} \left(\frac{256^2}{\frac{1}{ij} \sum_i \sum_j (I[i][j] - O[i][j])^2} \right).$$

Figure 7.5 shows the quality of the denoising for a set of Lena images with increasing quantity of noise. The curves are obtained using the median (M) and different values of k' in the discrepancy while k is fixed at 25. The median is better when the noise ratio is small but as we increase the number of outliers, the discrepancy obtains better results. This also shows that the optimal k' depends on the noise ratio. It also depends on the image we consider and thus makes it difficult to find an easy way to choose it automatically. Heuristically, it is better to take k' around $\frac{2}{3}k$, especially when there is a lot of noise.

State of the art results in computer vision obtain better experimental results (e.g. [54, 79, 96]). However, these results assume that the noise model is known and they can start by detecting and removing noisy points before rebuilding the image. Our methods are free from assumptions on the generative model of the image. The algorithms do not change depending

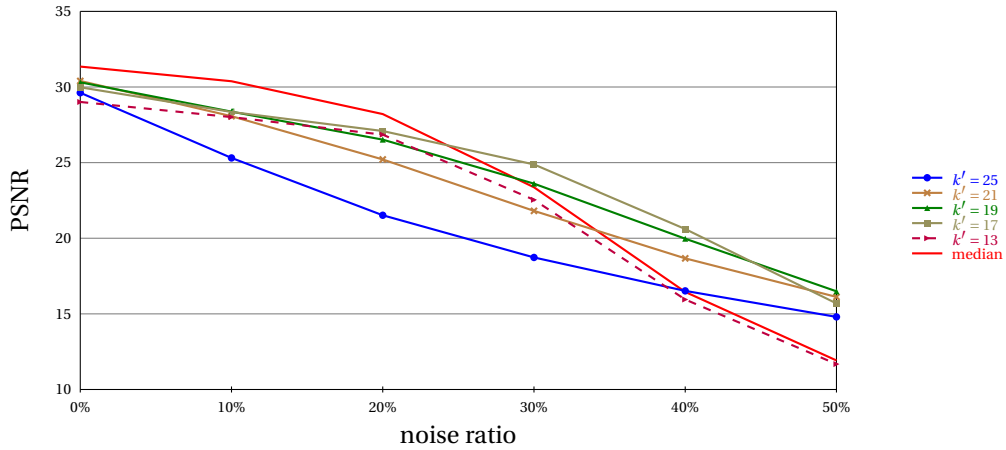


Figure 7.5 – PSNR for Lena images depending on the choice of k' and the quantity of noise

on the type of noise.

7.3 Scalar field analysis with unbounded geometric noise

For the rest of the chapter, we assume that M is a manifold with curvature bounded by c_M and positive reach r_M . Moreover, we assume that the input P is an (ε, r) -sample of M , where

$$\varepsilon \leq \frac{r_M}{6} \text{ and } r > 2\varepsilon \quad (7.3)$$

We assume that there is no intrinsic functional noise in the sense that for any $p \in P$, the observed function value $\tilde{f}(p) = f(\pi(p))$ is the same as the true value for the projection $\pi(p) \in M$ of this point. Remark that $\pi(p)$ is not always well-defined. If there is an ambiguity, we choose arbitrarily among the possible projections.

Taking advantage of the interleaving from Lemma 6.2, we can use the distance to the empirical measure to filter the points of P . In particular, we consider the set

$$L = P \cap d_{\mu, m}^{-1}([-\infty, \eta]) \text{ where } \eta \geq 2\varepsilon.$$

We then use a similar approach as the one from [32] for this set L . The optimal choice for the parameter η is 2ε . However, any value with $\eta \leq r$ and $\eta + \varepsilon < \rho_M$ works as long as there exist δ and δ' satisfying the conditions stated in Theorem 7.3.

Let $\tilde{L} = \{\pi(x) | x \in L\}$ denote the orthogonal projection of L onto M . To simulate sub-level sets, consider the restricted sets $L_\alpha := L \cap (f \circ \pi)^{-1}([-\infty, \alpha])$ and let $\tilde{L}_\alpha = \pi(L_\alpha)$. By our assumption on \tilde{f} , we have: $L_\alpha = \{x \in L | \tilde{f}(x) \leq \alpha\}$.

Let us recall a result about the relation between Riemannian and Euclidian metrics [50]. For any two points $x, y \in M$ with $d(x, y) \leq \frac{r_M}{2}$,

$$d(x, y) \leq d_M(x, y) \leq \left(1 + \frac{4d(x, y)^2}{3r_M^2}\right) d(x, y) \leq \frac{4}{3} d(x, y). \quad (7.4)$$

As a direct consequence of the noise model, for any point $x \in M$, there exists a point $y \in L$ at distance less than 2ϵ . We will use the *extrinsic* Vietoris-Rips complex to infer the scalar field topology. Using the previous relation, we obtain that for points in L , the Euclidean distance for nearby points approximates the Riemannian metric on M .

Proposition 7.11 *Let $\lambda = \frac{4}{3} \frac{r_M}{r_M - (\eta + \epsilon)}$ with $\eta + \epsilon \leq r_M$. For $\forall (x, y) \in L \times L$, $d(x, y) \leq \frac{r_M}{2} - \frac{\eta + \epsilon}{2}$ implies*

$$\frac{d_M(\pi(y), \pi(x))}{\lambda} \leq d(x, y) \leq 2(\eta + \epsilon) + d_M(\pi(x), \pi(y)).$$

Proof: Let x and y be two points of L such that $d(x, y) \leq \frac{r_M}{2} - \frac{\eta + \epsilon}{2}$. As $d_{\mu, m}(x) \leq \eta \leq r$, and μ is an (ϵ, r) -sample of M , $d(\pi(x), x) \leq \eta + \epsilon$. Therefore $d(\pi(x), \pi(y)) \leq \frac{r_M}{r_M - (\eta + \epsilon)} d(x, y)$ [61, Theorem 4.8, (8)]. This implies $d(\pi(x), \pi(y)) \leq \frac{r_M}{2}$ and following (7.4), $d_M(\pi(x), \pi(y)) \leq \frac{4}{3} d(\pi(x), \pi(y))$. ■

Theorem 7.12 *Let M, f, L be defined as previously. Then, for any $\delta \geq 2\eta + 6\epsilon$ and any $\delta' \in \left[2\eta + 2\epsilon + \frac{8}{3} \frac{r_M}{r_M - (\eta + \epsilon)} \delta, \frac{3}{4} \frac{r_M - (\eta + \epsilon)}{r_M} \varrho(M)\right]$, $H_*(f)$ and $H_*(R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha))$ are $\frac{4}{3} \frac{cr_M \delta'}{r_M - (\eta + \epsilon)}$ -interleaved.*

Proof: First, note that \bar{L} is a 2ϵ -sample of M in its Riemannian metric. This is because for any point $x \in M$, we know that there exists some $p \in L$ such that $d(x, p) \leq d_{\mu, m}(x) \leq \epsilon$. Hence $d(x, \pi(p)) \leq d(x, p) + d(p, \pi(x)) \leq 2d(x, p) \leq 2\epsilon$. Now we apply Theorem 7.3 to \bar{L} by using $\tilde{d}(\pi(x), \pi(y)) := d(x, y)$; and setting $\lambda = \mu = \frac{4}{3} \frac{r_M}{r_M - (\eta + \epsilon)}$, $\nu = 2(\eta + \epsilon)$: the requirement on the distance function \tilde{d} in Theorem 7.3 is satisfied due to Proposition 7.11. ■

Since M is compact, f is bounded due to the Lipschitz condition. We can look at the limit when $\alpha \rightarrow \infty$. There exists a value T such that for any $\alpha \geq T$, $L_\alpha = L$ and $f^{-1}([-\infty, \alpha]) = M$. The above interleaving means that $H_*(M)$ and $H_*(R_\delta(L) \hookrightarrow R_{\delta'}(L))$ are interleaved. However, both objects do not depend on α and this gives the following inference result:

Corollary 7.13 *$H_*(M)$ and $H_*(R_\delta(L) \hookrightarrow R_{\delta'}(L))$ are isomorphic.*

7.4 Scalar field analysis with both functional and geometric noise

Our constructions can be combined to analyze scalar fields in a more realistic setting. Assuming that the point set P is an (ϵ, r) -sampling of M , we adapt the condition of Definition 7.4 to take into account the geometry and we assume that there exist η and s such that:

$$\forall p \in d_{\mu, m}^{-1}([-\infty, \eta, 1]), |\{q \in NN_k(p) \mid |\tilde{f}(q) - f(\pi(p))| \leq s\}| \geq k' \quad (7.5)$$

Note that we are using $f(\pi(p))$ as the “true” function value at a sample p which is off the manifold M . The condition on the functional noise is only for points close to the manifold and hence with small distance to measure. Combining the methods from the previous two sections, we obtain the *combined noise algorithm* where η is a parameter greater than ϵ .

We propose the following 3-steps algorithm. It starts by handling outliers in the geometry then it makes a regression on the function values before running the existing algorithm for scalar field analysis [32] on the filtration $\hat{L}_\alpha = \{p \in L \mid \hat{f}(p) \leq \alpha\}$.

COMBINED NOISE ALGORITHM

1. Compute $L = P \cap d_{\mu, m}^{-1}([-\infty, \eta])$.
2. Replace functional values \tilde{f} by \hat{f} using either k -median or discrepancy based method.
3. Run the scalar field analysis algorithm from [32] on (L, \hat{f}) .

Theorem 7.14 *Let M be a manifold, f a c -Lipschitz function on M . Let P be an (ϵ, r) -sampling of M and \tilde{f} functional values such that (7.5) is satisfied, where $\eta \geq \epsilon$. The combined noise algorithm has the following guarantees.*

For any $\delta \in [2\eta + 6\epsilon, \frac{\varrho(M)}{2}]$ and any $\delta' \in [2\eta + 2\epsilon + \frac{8}{3} \frac{r_M}{r_M - (\eta + \epsilon)} \delta, \frac{3}{4} \frac{r_M - (\eta + \epsilon)}{r_M} \varrho(M)]$, $H_(f)$ and $H_*(R_\delta(\hat{L}_\alpha) \hookrightarrow R_{\delta'}(\hat{L}_\alpha))$ are $(\frac{4}{3} \frac{cr_M \delta'}{r_M - (\eta + \epsilon)} + \xi s)$ -interleaved where $\xi = 1$ if we use the median and $\xi = (1 + 2\sqrt{\frac{k-k'}{2k'-k}})$ if we use the discrepancy.*

Notice that while this theorem assumes a setting where we can ensure theoretical guarantees, the algorithm can be applied in a more general setting and still produce good results. Beware that Theorems 7.10 and 7.12 do not combine directly. Applying either of them will make it impossible to guarantee that the hypotheses of the other one are still verified. It is necessary to directly combine parts of the proofs.

Proof: First, consider the two filtration $\{F_\alpha\}$ and $\{L_\alpha\}$ where $L_\alpha = \{x \in L \mid f(\pi(x)) \leq \alpha\}$. Using Theorem 7.12, for any $\delta \geq \nu + 2\mu \frac{\beta}{\lambda}$ and any $\delta' \in [2\eta + 2\epsilon + \frac{8}{3} \frac{r_M}{r_M - (\eta + \epsilon)} \delta, \frac{3}{4} \frac{r_M - (\eta + \epsilon)}{r_M} \varrho(M)]$, $H_*(f)$ and $H_*(R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha))$ are $\frac{4}{3} \frac{cr_M \delta'}{r_M - (\eta + \epsilon)}$ -interleaved.

Consider $\hat{L}_\alpha = \{p \in L \mid \hat{f}(p) \leq \alpha\}$. Our algorithm returns $H_*(R_\delta(\hat{L}_\alpha) \hookrightarrow R_{\delta'}(\hat{L}_\alpha))$. Fix α and let (x, y) be an edge of $R_\delta(L_\alpha)$. Then $d(x, y) \leq 2\delta$, $f(\pi(x)) \leq \alpha$, $f(\pi(y)) \leq \alpha$. Corollary 7.9 can be applied to $f \circ \pi$ due to hypothesis (7.5). Hence $|\hat{f}(x) - f(\pi(x))| \leq \xi s$ and $|\hat{f}(y) - f(\pi(y))| \leq \xi s$. Thus $(x, y) \in R_\delta(\hat{L}_{\alpha + \xi s})$. One can reverse the role of \hat{f} and f and get an ξs -interleaving of $\{R_\delta(L_\alpha)\}$ and $\{R_\delta(\hat{L}_\alpha)\}$. We have two filtrations of the same metric space that are interleaved. All parts of the diagram in Figure 7.6 commute as all the arrows are induced by inclusions.

Thus the nested filtrations are interleaved and we obtain that $H_*(R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha))$ and $H_*(R_\delta(\hat{L}_\alpha) \hookrightarrow R_{\delta'}(\hat{L}_\alpha))$ are ξs -interleaved. ■

We now provide a practical example as a proof of concept. We consider an elevation map of an area near Corte in the French island of Corsica. The true measures of elevation are given in Figure 7.7. The topography can be analysed by looking at the function minus-altitude. We add random faulty sensors that give false results with a 20% probability to simulate malfunctioning equipment. The area covers a square of 2 minutes of arc in both latitude and longitude. We use the algorithm taking $k = 9$, $k' = 7$, $\eta = .05$ minute and $\delta = .025$ minute. We show the recovered diagrams in Figure 7.8, where the prominent peaks of the original image are highlighted. The gap is the ratio between the shortest living relevant feature and the longest feature created by the noise.

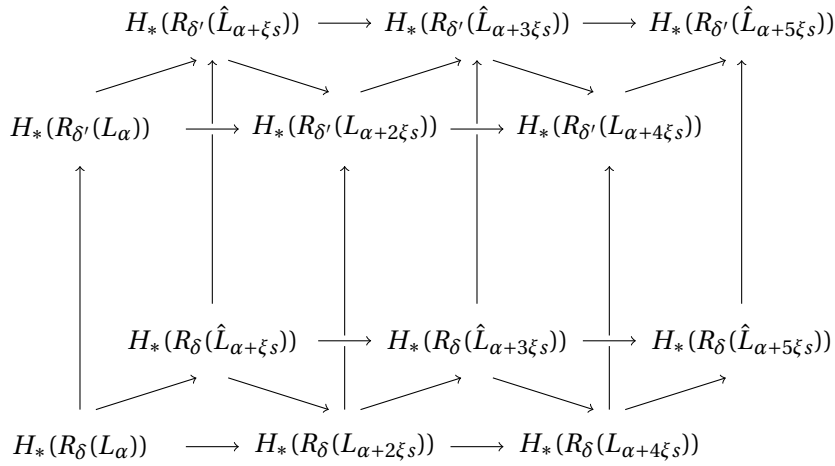


Figure 7.6 – Commutative diagram at homology level

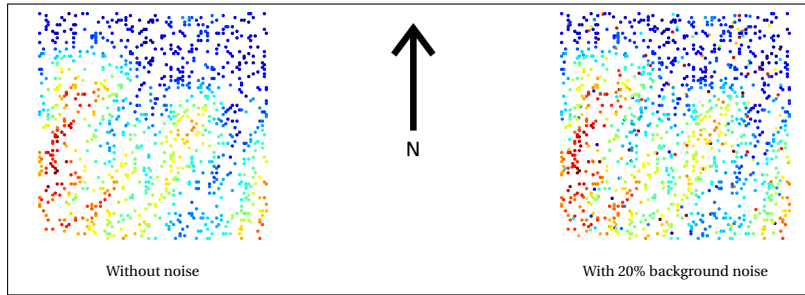


Figure 7.7 – Elevation map around Corte

Remark that the gap in the case of the noisy point cloud is less than 1. This means that one relevant topological feature has a shorter lifespan than one caused by noise. It is then impossible to recover the correct structure. The three methods are doing the scalar field analysis after a denoising of the functional values. In the case of the k -NN regression, the topological features are in the right order. However, the prominence given by the gap is significantly smaller than the one from the original point cloud. Both the discrepancy based method and the median provides gaps on par with the non-noisy input and thus allow a good recovery of the correct topology.

7.4. Scalar field analysis with both functional and geometric noise

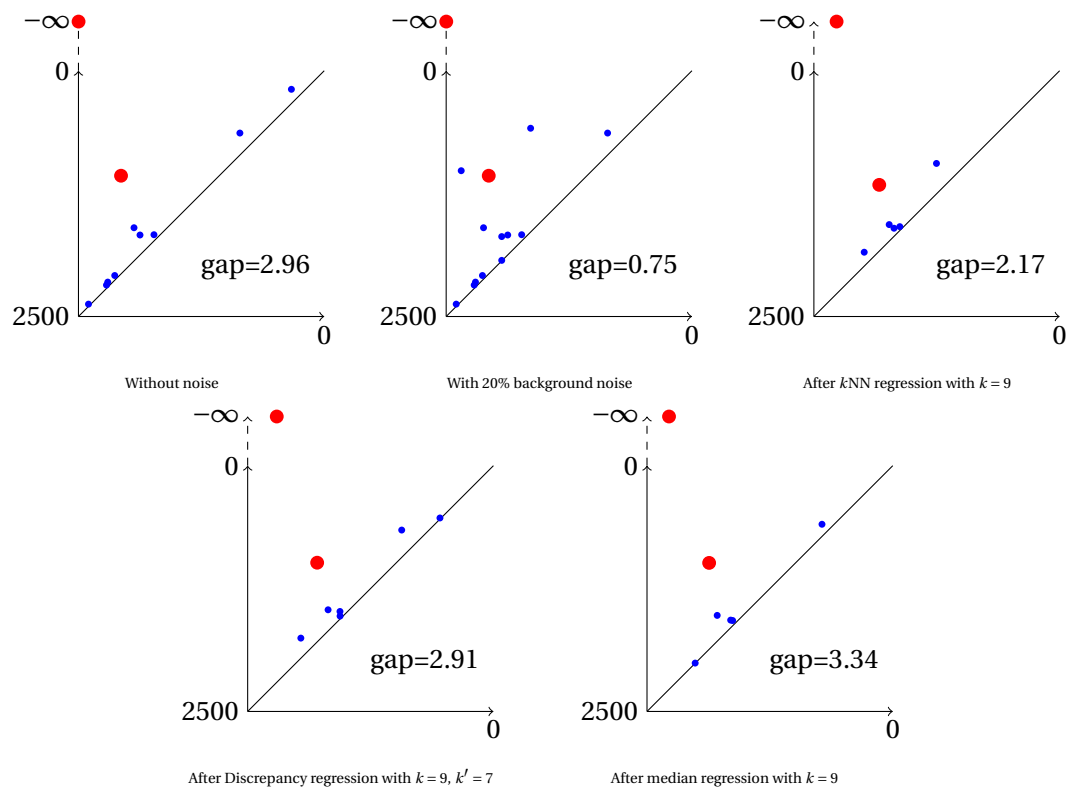


Figure 7.8 – Persistence diagrams of Corte Elevation map

8 Regression and incomplete data

The discrepancy introduced for the scalar field analysis can also be seen as a regression estimator. In this chapter, we explore the properties of this estimator from a statistical point of view and provide convergence rates for the discrepancy regression. Then, we introduce the problem of incomplete data as a direction for further research. We present an empirical way of using the regression estimator to tackle it.

8.1 Discrepancy is a regression estimator

First, we study the regression for the impulse noise model used in Section 7.2. Then, we discuss how it can be related to other noise models. Given a probability measure μ on \mathbb{R}^d and a bounded t -Lipschitz function $f : \text{Supp}(\mu) \rightarrow \mathbb{R}$, we consider $(X_i, Y_i)_{i \in [1, n]}$ independent and identically distributed variables obtained according to the law μ for X_i and conditionally to X_i :

$$Y_i = \begin{cases} f(X_i) & \text{with probability } 1 - \rho \\ Z_i & \text{with probability } \rho \end{cases}$$

where $\rho < \frac{1}{2}$ and Z_i has a distribution ν_x conditionally to $X_i = x$. Moreover, we assume that there exists C such that for all $x \in \text{Supp}(\mu)$, $\text{Supp}(\nu_x) \subset [-C, C]$. We introduce $K = \max(\|C - f\|_\infty, \max_{x, x'} |f(x) - f(x')|) < \infty$.

8.1.1 The discrepancy regression

We use the discrepancy based estimator defined in Section 7.2.2 to obtain a regression estimator. We write \mathcal{B}_j the set of permutations of $[1, j]$. Given a point x , we index the X_i such that $X_{(i, n)}$ is the i^{th} nearest neighbour of x when we have picked n couples (X_i, Y_i) .

Definition 8.1 *Given an integer $k_n \leq n$ and $\alpha \in]0, 1]$, the discrepancy regression estimator is defined as*

$$f_n(x) = \operatorname{argmin}_{z \in \mathbb{R}} \min_{\sigma \in \mathcal{B}_{k_n}} \frac{1}{\lceil \alpha k_n \rceil} \sum_{i=1}^{\lceil \alpha k_n \rceil} \|Y_{(\sigma(i), n)} - z\|^2$$

Our aim is to recover the function f . Remark that this is not the classical case of a regression mixture where we want to recover $\mathbb{E}[Y|X]$, which is, writing $\bar{\nu}_x$ the mean of ν_x ,

$$\mathbb{E}[Y|X = x] = (1 - \rho)f(x) + \rho\bar{\nu}(x).$$

We consider the asymptotic behaviour of the estimator when:

$$\begin{cases} n & \rightarrow \infty \\ k_n & \rightarrow \infty \\ \frac{k_n}{n} & \rightarrow 0 \end{cases}$$

We separate points in two groups, the *good points* and the *outliers*.

Definition 8.2 Given $s \in \mathbb{R}$, $k_n \leq n$ and $x \in \mathbb{X}$, the set of good points among the k_n nearest neighbours of x is:

$$G_{(k_n, n)}^s(x) = \{X_{(i, n)}(x) | i \leq k_n \text{ and } |Y_{(i, n)} - f(X_{(i, n)})| \leq s\}$$

Points that are not good for all $x \in \mathbb{X}$ are called outliers. We bound the probability for any x to have more than $(1 - \alpha)k_n$ outliers among its k_n nearest neighbours.

Lemma 8.3 Given $x \in \mathbb{X}$ and $\frac{1}{2} < \alpha < 1 - \rho$,

$$P_B = \mathbb{P} \left[\frac{|G_{(k_n, n)}^0|}{k_n} < \alpha \right] < e^{-2(1-\alpha-\rho)^2 k_n}$$

Proof:

$$P_B = \mathbb{P} \left[\frac{|G_{(k_n, n)}^0|}{k_n} < \alpha \right] = \mathbb{P} \left[1 - \rho - \frac{|G_{(k_n, n)}^0|}{k_n} > 1 - \rho - \alpha \right]$$

Hence, using Hoeffding inequality for a Bernoulli law, $P_B < e^{-2(1-\alpha-\rho)^2 k_n}$. ■

8.1.2 Convergence rate for discrepancy

Restricting ourselves to the case where the dimension d of the ambient space is greater than 3, we introduce [68, Lemma 6.4].

Lemma 8.4 Assume that X is bounded and $d \geq 3$ then there exists a constant \tilde{c} such that:

$$\mathbb{E} [||X_{(1, n)}(X) - X||^2] \leq \frac{\tilde{c}}{n^{\frac{2}{d}}}$$

In our setting, we obtain the convergence rate:

Theorem 8.5 For any $x \in \text{Supp}(\mu)$ and $\frac{1}{2} < \alpha < 1 - \rho$,

$$\mathbb{E}(f_n(x) - f(x))^2 \leq e^{-2(\alpha+\rho-1)^2 k_n} K^2 + \left(1 + \sqrt{\frac{1-\alpha}{2\alpha-1}}\right)^2 \tilde{c} t^2 \left(\frac{1+2k_n}{n}\right)^{\frac{2}{d}}$$

Proof: Let x be a point of $\text{Supp}(\mu)$. We will split occurrences in two categories. Either we have more than $(1 - \alpha)k_n$ outliers in the k_n nearest neighbours of x or not. If we have too many outliers, we simply bound $|f_n(x) - f(x)|$ by the constant K . In the other case, we locally have a

$(k_n, \alpha k_n, t||X_{(k_n, n)} - x||)$ -functional-sample of f and apply Lemma 7.8.

$$\mathbb{E}(f_n(x) - f(x))^2 \leq P_B K^2 + (1 - P_B) \left(1 + \sqrt{\frac{1 - \alpha}{2\alpha - 1}}\right)^2 t^2 \mathbb{E}[||X_{(k_n, n)} - x||^2] = I_1(x) + I_2(x)$$

Now split the data X_1, \dots, X_n into $2k_n + 1$ segments such that the first $2k_n$ segments have length $\lfloor \frac{n}{2k_n} \rfloor$ and let \tilde{X}_j^x be the nearest neighbour of x from the j^{th} segment. Then $\tilde{X}_1^x, \dots, \tilde{X}_{2k_n}^x$ are $2k_n$ different elements of $\{X_1, \dots, X_n\}$ which implies:

$$||X_{(k_n, n)} - n|| \leq \frac{1}{2k_n} \sum_{i=1}^{2k_n} ||X_{(i, n)} - x|| \leq \frac{1}{2k_n} \sum_{i=1}^{2k_n} ||\tilde{X}_i^x - x||$$

Noticing that $\frac{1}{l} (\sum_{i=1}^l a_i)^2 \leq \sum_{i=1}^l a_i^2$ and using Jensen inequality:

$$\begin{aligned} I_2(x) &\leq \left(1 + \sqrt{\frac{1 - \alpha}{2\alpha - 1}}\right)^2 t^2 \mathbb{E} \left[\left(\frac{1}{2k_n} \sum_{j=1}^{2k_n} ||\tilde{X}_j^x - x|| \right)^2 \right] \\ &\leq \left(1 + \sqrt{\frac{1 - \alpha}{2\alpha - 1}}\right)^2 \frac{t^2}{2k_n} \sum_{j=1}^{2k_n} \mathbb{E} [||\tilde{X}_j^x - x||^2] \\ &= \left(1 + \sqrt{\frac{1 - \alpha}{2\alpha - 1}}\right)^2 t^2 \mathbb{E} [||\tilde{X}_1^x - x||^2] \\ &= \left(1 + \sqrt{\frac{1 - \alpha}{2\alpha - 1}}\right)^2 t^2 \mathbb{E} [||\tilde{X}_{(1, \lfloor \frac{n}{2k_n} \rfloor)} - x||^2] \end{aligned}$$

Using Lemma 8.4,

$$I_2(x) \leq \left(1 + \sqrt{\frac{1 - \alpha}{2\alpha - 1}}\right)^2 \frac{t^2 \tilde{c}}{\left[\frac{n}{2k_n}\right]^{\frac{2}{d}}}$$

■

8.1.3 Convergence rate for median

We can also define a regression estimator using the median.

Definition 8.6 Given an integer $k_n \leq n$, the median regression is defined by $f_n^\Delta(x)$, the median of the values $(Y_{(i, n)})$ for $i \in [1, k_n]$.

Theorem 8.7 For any $x \in \text{Supp}(\mu)$,

$$\mathbb{E}(f_n^\Delta(x) - f(x))^2 \leq e^{-2(\frac{1}{2} - \rho)^2 k_n} K^2 + \tilde{c} t^2 \left(\frac{1 + 2k_n}{n} \right)^{\frac{2}{d}}$$

Proof: We apply the same kind of proof as for the discrepancy convergence rates. We separate the cases whether at least half the points are good or not. The probability to have more than half the points that are outliers is given by the probability P_H obtained by applying Hoeffding

inequality:

$$P_H = \mathbb{P} \left[\frac{|G_{(k_n, n)}^s|}{k_n} < \frac{1}{2} \right] \leq \mathbb{P} \left[1 - \rho - \frac{|G_{(k_n, n)}^s|}{k_n} > \frac{1}{2} - \rho \right] \leq e^{-2(\frac{1}{2} - \rho)^2 k_n}$$

In the case where at least half the points are good, we can use the Lipschitz condition on f to bound the error. It will be at most $t\|x - X_{(k_n, n)}\|$.

$$\mathbb{E}(f_n^\Delta(x) - f(x))^2 \leq P_H K^2 + (1 - P_H) t^2 \mathbb{E}\|x - X_{(k_n, n)}\|^2$$

The result is then obtained through Lemma 8.4 and the same proof as in Theorem 8.5. \blacksquare

8.1.4 Relaxing the noise model

Our noise model is not a classical one. It highlights the kind of problem that can be found with k -NN regression for example. We try to recover a value that is not the conditional expectation of Y knowing X and thus the k -NN regression can present a bias.

We consider more classical noise models. We do not change the probability law for X_i and $Y_i = f(X_i) + \epsilon_i$ where ϵ_i has probability distribution ν .

Theorem 8.8 *Given $\alpha \in]0, 1]$, if $\operatorname{argmin}_{\mathbb{R}}(d_{\nu, \alpha}) = \{0\}$ and $\int_{\mathbb{R}} |z|^3 \nu(dz) < \infty$, then for any $x \in \operatorname{Supp}(\mu)$ and $\epsilon > 0$,*

$$\lim_{k_n \rightarrow \infty; \frac{k_n}{n} \rightarrow 0} \mathbb{P} [|f_n(x) - f(x)| > \epsilon] = 0$$

Proof: Using Lemma 8.4, we know that

$$\mathbb{E} [\|\tilde{X}_i^x - x\|^2] \leq \frac{\tilde{C}}{\left\lfloor \frac{n}{k_n} \right\rfloor^{\frac{2}{d}}}.$$

We denote $\tilde{\nu}$ the translated probability measure $f(x) + \nu$, y_n the empirical measure $\frac{1}{k_n} \sum_{i=1}^{k_n} \tilde{Y}_i^x$ and $\hat{\nu}_n = y_n - f(x)$. Given that f is t -Lipschitz and thanks to Theorem 6.6:

$$\begin{aligned} \mathbb{E} [W_2(\tilde{\nu}, y_n)^2] &\leq 2\mathbb{E} \left[\frac{1}{k_n} \sum_{i=1}^{k_n} \|f(x) - \tilde{Y}_i^x\|^2 \right] + 2\mathbb{E}[W_2(\nu, \hat{\nu}_n)^2] \\ &\leq 2t^2 \mathbb{E} [\|\tilde{X}_i^x - x\|^2] + 2\mathbb{E}[W_2(\nu, \hat{\nu}_n)^2] \\ &\leq 2t^2 \frac{\tilde{C}}{\left\lfloor \frac{n}{k_n} \right\rfloor^{\frac{2}{d}}} + 2 \frac{C}{k_n^{\frac{2}{5}}} \end{aligned}$$

for some constant C . Hence $\mathbb{E} [\|d_{\tilde{\nu}, \alpha} - d_{y_n, \alpha}\|_\infty] \leq \frac{1}{\sqrt{\alpha}} \mathbb{E}[W_2(\tilde{\nu}, y_n)] \rightarrow 0$. Furthermore, $d_{\tilde{\nu}, \alpha}$ is proper as shown in the proof of Proposition 4.18 and $\operatorname{argmin}_{\mathbb{R}} d_{\tilde{\nu}, \alpha} = f(x)$. Therefore, there exists a sequence of positive real numbers $\{\epsilon_n\}_{n>0}$ and $N > 0$ such that $\lim_{n \rightarrow \infty} \epsilon_n = 0$ and for

any $n \geq N$, $d_{\tilde{v},\alpha}^{-1}([0, f(x) + \frac{1}{n}]) \subset [-\epsilon_n, \epsilon_n]$. Using Markov inequality:

$$\mathbb{P}\left[\|d_{\tilde{v},\alpha} - d_{y_n,\alpha}\|_\infty > \frac{\epsilon_n}{2}\right] \leq \frac{2\mathbb{E}[\|d_{\tilde{v},\alpha} - d_{y_n,\alpha}\|_\infty]}{\epsilon_n}$$

Given that $f_n(x) = \operatorname{argmin}_{\mathbb{R}} d_{y_n,\alpha}$ and $f(x) = \operatorname{argmin}_{\mathbb{R}} (d_{\tilde{v},\alpha})$,

$$\mathbb{P}[|f_n(x) - f(x)| > \epsilon_n] \leq \frac{2\mathbb{E}[\|d_{\tilde{v},\alpha} - d_{y_n,\alpha}\|_\infty]}{\epsilon_n}.$$

■

This can be in particular applied to the popular Gaussian noise model where ν is the probability measure of the normal law $\mathcal{N}(0, \sigma^2)$.

Corollary 8.9 *If ν is the probability measure of the normal law $\mathcal{N}(0, \sigma^2)$, then for any $x \in \operatorname{Supp}(\mu)$, $\epsilon > 0$ and $\alpha \in]0, 1]$,*

$$\lim_{k_n \rightarrow \infty; \frac{k_n}{n} \rightarrow 0} \mathbb{P}[|f_n(x) - f(x)| > \epsilon] = 0$$

Proof: We show that $\operatorname{argmin}(d_{\nu,m}) = 0$ for any mass $m > 0$ and is uniquely defined. In fact, for any radius $r > 0$ and $x \neq 0$, $\nu(B(x, r)) < \nu(B(0, r))$. Moreover, ν is a measure with density and thus $\delta_{\nu,l}(0) > \delta_{\nu,l}(x)$, because $\nu(B(x, r + \epsilon)) \rightarrow \nu(B(x, r))$ when $\epsilon \rightarrow 0$. Hence, $d_{\nu,m}(0) < d_{\nu,m}(x)$. Therefore $\operatorname{argmin}(d_{\nu,m})$ and is uniquely defined and is zero.

Moreover the third moment of $\mathcal{N}(0, r)$ is bounded and Theorem 8.8 applies. ■

However, $\operatorname{argmin}_{\mathbb{R}} d_{\nu,\alpha}$ is not always uniquely defined for some usual noise models, the tubular noise model for example. In this case, we can guarantee that $f_n(x)$ is a good approximation of $f(x)$.

Theorem 8.10 *If ν is the uniform probability measure on $[-\delta, \delta]$ for some $\delta > 0$, then for any $x \in \operatorname{Supp}(\mu)$ and $\alpha \in]0, 1]$,*

$$\mathbb{E}[(f_n(x) - f(x))^2] \leq \left(1 + \sqrt{\frac{1-\alpha}{2\alpha-1}}\right)^2 \left(\tilde{c} t^2 \left(\frac{1+2k_n}{n}\right)^{\frac{2}{d}} + \delta^2\right)$$

Proof: The proof is identical to the one of Theorem 8.5 when $P_B = 0$ and replacing the local $(k_n, \alpha k_n, t||X_{(k_n,n)} - x||)$ -functional-sample of f by a local $(k_n, \alpha k_n, t||X_{(k_n,n)} - x|| + \delta)$ -functional-sample. The added δ comes from the uncertainty created by the tubular noise. ■

If the support of μ is embedded in a Riemannian manifold, we can adapt the proofs as long as there exists a property similar to Lemma 8.4. This is the case for Riemannian submanifold of \mathbb{R}^d . Moreover, if we have Riemannian submanifold of \mathbb{R}^d with positive reach and bounded curvature, we can use Proposition 7.11 to guarantee that the regression estimator will converge even if we only know the extrinsic metric.

8.2 Application to incomplete data

In this section, we propose an example where the data is incomplete. Consider a set of points in high dimension where some of the coordinates are missing. The missing coordinates are not the same for every point and we hope to recover the value of those missing coordinates and the structure of the underlying data by using the discrepancy regression.

This situation happens in real cases. During a poll, some people can forget to answer certain question or do not wish to answer them. Assuming that the troublesome questions are not the same for everyone, we obtain a set with incomplete data. Similarly, if we try to help for medical diagnosis, not every patient will have undergone the whole set of available tests. Each patient is represented as a point whose coordinates are the results of medical tests. Some of the patients miss certain tests. Therefore, we have an incomplete data case.

This section aims to illustrate how tools used in the scalar field analysis of Chapter 7 can be adapted to the setting of incomplete data.

8.2.1 Algorithm for recovery of incomplete data

We showed in Section 7.2.3 how the discrepancy-based method could be used to recover images with aberrant values. Missing coordinates can be seen directly as aberrant values and, for example, be replaced by a random value before applying the scalar field algorithm. However, we want to try to take advantage of knowing which coordinates are missing.

Consider a set of points $P \in \mathbb{R}^d$. We assume that there exists a t -Lipschitz function $f: \mathbb{R}^{d'} \rightarrow \mathbb{R}^d$ and a point set $Q \in \mathbb{R}^{d'}$ such that $P = f(Q)$. The dimension d' can be smaller than d and is the intrinsic dimension of our point cloud.

Our algorithm starts by normalising all dimensions in order to give them the same weight. Then, we introduce a notion of distance between two points $x = (x_1, \dots, x_d)$ and $y = (y_1, \dots, y_d)$. The distance is the average of the sum of squares for the known coordinates common to x and y . If the coordinates $x_{i_1}, x_{i_2}, y_{i_1}$ and y_{i_3} are missing coordinates. Then the distance $\tilde{d}(x, y)$ between x and y is defined as:

$$\tilde{d}(x, y) = \sqrt{\frac{1}{d-3} \sum_{i \in [1, d] \setminus \{i_1, i_2, i_3\}} |x_i - y_i|^2}$$

Given this notion of distance, we use the discrepancy with two parameters k and k' to infer a value for the missing coordinate. Consider a missing coordinate x_i of the point x . Among the k -nearest neighbours of x , we consider the set of points whose i^{th} coordinates is known. If the cardinality of this set is more than k' , then we apply the discrepancy method. Otherwise, we take the mean of these values.

The algorithm is elementary and tries to take into account the information at our disposal. It is able to handle some noise. Remark that in some cases, some missing values can not be recovered. We kept this algorithm simple and it should be taken as a first step into an extension of the methods from Chapter 7.

MISSING COORDINATES RECOVERY

1. Normalise the point cloud P .
 2. Compute the distance matrix of P .
 3. For every missing value x_i :
 - (a) Compute the set of points among the k -nearest neighbours of x whose i^{th} is known.
 - (b) If this set empty, indicate that we cannot recover the value, otherwise compute the new value.
-

8.2.2 Illustration on a synthetic example

We give some experimental results on a synthetic example. We consider a point set in \mathbb{R}^{10} with intrinsic dimension 3. It is obtained by sampling uniformly the square $[-10, 10]^3$ and applying the function $f : \mathbb{R}^3 \mapsto \mathbb{R}^{10}$ defined for any $q = (x, y, z)$ by,

$$\begin{aligned}
 f(q) &= (f_1(q), \dots, f_{10}(q)) \\
 f_1(q) &= |x + y - 2z| \\
 f_2(q) &= \sqrt{|x^2 - y^2|} \\
 f_3(q) &= |x - z| \\
 f_4(q) &= 5x + 3y + 2z \\
 f_5(q) &= \cos(x) + \sin(y) \\
 f_6(q) &= e^{\frac{-(x-y)^2}{100}} \\
 f_7(q) &= \ln(1 + |z - x|) \\
 f_8(q) &= Re(\sqrt{1 + \frac{x + y - x^2}{4}}) \\
 f_9(q) &= (x^7 + y^7 + z^7)^{\frac{1}{7}} \\
 f_{10}(q) &= x \sin(z) - y \sin(z)
 \end{aligned}$$

We ran our algorithm on input datasets of 10000 points and using $k = 50$. The graphics of Figure 8.1 show the relative error on the recovered coordinates for various choices of the ratio $\frac{k'}{k}$. Missing coordinates are obtained by removing each coordinate with probability p between .01 and .35. In the noisy case, we added impulse noise. It means that we replaced 25 percent of the initial coordinates by a random value uniformly picked inside $[-1, 1]$ after normalisation. We recover the missing data with a relative error of 10 percent when the proportion of missing coordinates is less than 15 percent. The quality decreases with more missing data. Note that,

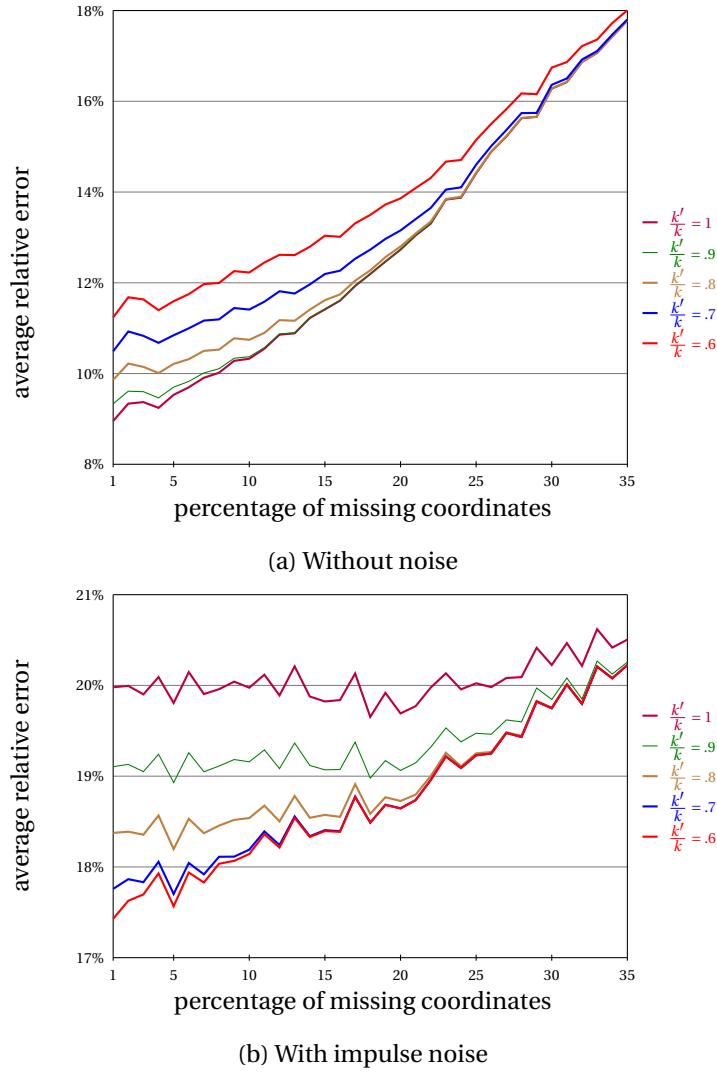


Figure 8.1 – Average relative error on recovered coordinates

in presence of noise, the larger the ratio $\frac{k'}{k}$, the more important is the influence of the noise. The missing data start to influence notably the quality of the recovery when the proportion of missing data is around the same as the proportion of noisy points for ratios $\frac{k'}{k}$ of .9 or 1. These experiments give an insight for the possible use of distance based methods, such as the discrepancy, for the treatment of data with missing coordinates. A more theoretical approach would need a more precise model for incomplete data and the way it influences the computation of the distance matrix.

9 Conclusion

In this thesis, we provided a complete method to approximate the persistence diagram of the distance to a measure with tractable complexity. Thus, it is now possible to use persistent homology in practice with data containing outliers and theoretical guarantees on the results. Furthermore, we introduced a new regression operator using the distance to a measure in order to analyse scalar fields and incomplete data. Some interesting questions are still unanswered.

Identifiability of measures We showed results on the identifiability of measures from distances to measures. The amount of information required is quite large. It seems that the reconstruction algorithms from a finite number of points could be adapted into a proof of the identifiability of a measure μ in \mathbb{R}^d knowing the distance to μ at $d + 2$ points. Moreover, the knowledge of the distance to a measure on the whole space for only one mass parameter might be sufficient to identify the measure.

These possible results are still qualitative. There exists no constant c such that $\|d_{\mu,m} - d_{\nu,m}\|_{\infty} \geq cW_2(d_{\mu,m}, d_{\nu,m})$ in all generality. Such a result would have a huge impact on the use of the distance to a measure for signature purposes. The way to do this restriction is not obvious and would help to better understand the relation between distances to measures and Wasserstein distances.

Approximation of the distance to a measure Our work relies on power distances to approximate distances to measures. Lower bounds for distance to a measure approximation by power distances [80] are not reached by our approximations. A deterministic algorithm to attain an arbitrarily additive approximation of the distance to a measure could help reach the lower bounds and offer a useful control mechanism. Possible tools to build such an algorithm include furthest point sampling and incrementally introducing points using the structure of the measure.

Incomplete data framework Analysing incomplete data, as presented in Chapter 8, is very interesting for applications. No theoretical framework currently exists out of a linear setting to properly define the problem. We need to build a reasonable model for incomplete data before we can produce theoretical guarantees. Remark that our approach uses a very naive way of

computing the distance between two points. Improvements in this direction are without doubt possible.

A slight variation of this problem which may be easier to solve concerns missing distances. Forgetting about the coordinates, we can consider a set of points and the distance matrix between those points. Persistent homology can be used without knowing the coordinates, just the distances. Assume now that some of these distances are unknown. The question of how to do topological data analysis on this partial distance matrix is open. Again, it will need a well defined model of how the distances are missing from the matrix.

Statistical behaviour of persistent homology Considering persistence landscapes [16] instead of persistence diagrams makes it possible to look at persistent homology from a statistical point of view. Recent work uses this approach to derive convergence rates [30, 31] and confidence bounds [60]. It is also possible to design subsampling methods [28] or use bootstrapping [29]. These results are a first step towards the elaboration of a statistical framework for persistence diagrams and landscapes, but much remains to be done in order to have usable statistical models about persistent homology.

A Proof of Theorem 5.12

Proof: Let us fix three distinct points a, b and c . We will show that whatever weights are put on these three points, the three relations of the triangle inequality hold.

General remarks: First, we consider some general remarks used through the proof. Due to the metric structure of \mathbb{X} , we have the triangle inequality:

$$d_{\mathbb{X}}(a, b) \leq d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c) \quad (\text{A.1})$$

Remark also that the function f is never less than half the distance in the metric space \mathbb{X} . Intuitively, this is a consequence of the definition as an intersection of two balls. Even with zero weights, the union of ball has to cover the space between the points and thus the sum of their radii need to be at least equal to the distance between the two points. Due to our definition of the power distance, the parameter α is always smaller than the radii. This translates in more formal language, for any $(a, b) \in P^2$,

$$f(a, b) \geq \frac{d_{\mathbb{X}}(a, b)}{2} \quad (\text{A.2})$$

Let us assume that $|w_a^2 - w_b^2| < d_{\mathbb{X}}(a, b)^2$. Then, by definition:

$$f(a, b) \geq \sqrt{\frac{d_{\mathbb{X}}(a, b)^2}{4}} = \frac{d_{\mathbb{X}}(a, b)}{2}.$$

If our assumption is not correct, then, choosing a and b such that $w_b \geq w_a$, we get:

$$f(a, b) = w_b \geq \sqrt{w_b^2 - w_a^2} \geq d_{\mathbb{X}}(a, b)^2.$$

We will now explore the different cases for the expression of f . We say that a pair of points (a, b) is saturated if $|w_a^2 - w_b^2| \geq d_{\mathbb{X}}(a, b)^2$. This implies that the distance $f(a, b) = \max(w_a, w_b)$.

All saturated First, we consider the case where all pairs of points are saturated. Without loss of generality, we assume that $w_c \geq w_b \geq w_a$. Then, $f(a, b) = w_b$ and $f(a, c) = f(b, c) = w_c$.

The triangle inequality for f is given by:

$$\begin{aligned} f(a, b) &\leq f(a, c) \leq f(a, c) + f(b, c) \\ f(b, c) &= f(a, c) \leq f(a, c) + f(a, b) \\ f(a, c) &= f(b, c) \leq f(b, c) + f(a, b) \end{aligned}$$

All saturated but one We now release one of the pair. Without loss of generality, we assume that c is the common point to the two saturated pairs and that $w_b \geq w_a$. This gives the three following relations:

$$w_b^2 - w_a^2 < d_{\mathbb{X}}(a, b)^2 \quad (\text{A.3})$$

$$|w_c^2 - w_a^2| \geq d_{\mathbb{X}}(a, c)^2 \quad (\text{A.4})$$

$$|w_c^2 - w_b^2| \geq d_{\mathbb{X}}(b, c)^2 \quad (\text{A.5})$$

To complete the proof, we need to consider the three possible cases for the ordering of w_a , w_b and w_c .

Case 1.1: $w_c \geq w_b \geq w_a$

Let us remark that $f(a, c) = f(b, c)$. Then, we only need to check that $f(a, b) \leq 2f(a, c) = 2w_c$.

$$\begin{aligned} f(a, b)^2 &= \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} \\ &\leq \frac{w_c^2}{2} + \frac{w_c^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} && w_c \geq w_b \geq w_a \\ &\leq w_c^2 + \frac{2}{4}(d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2) + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} && (\text{A.1}); (x + y)^2 \leq 2(x^2 + y^2) \\ &\leq w_c^2 + \frac{d_{\mathbb{X}}(a, c)^2}{2} + \frac{d_{\mathbb{X}}(b, c)^2}{2} + \frac{w_b^2 - w_a^2}{4} && (\text{A.3}) \\ &\leq w_c^2 + \frac{w_c^2}{2} + \frac{w_c^2}{2} + \frac{w_b^2 - w_a^2}{4} && (\text{A.4}); (\text{A.5}) \\ &\leq 2w_c^2 + \frac{w_c^2}{4} && w_b \leq w_c \\ &\leq 4w_c^2 \end{aligned}$$

Case 1.2: $w_b \geq w_a \geq w_c$

Given the hypothesis, we have the relations:

$$f(a, c) = w_a$$

$$f(b, c) = w_b$$

Immediately we can deduce the first inequality:

$$f(a, c) \leq f(b, c) \leq f(b, c) + f(a, b)$$

Let us consider $f(a, b)$:

$$\begin{aligned} f(a, b)^2 &= \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} \\ &\leq \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2}{4} + \frac{w_b^2 - w_a^2}{4} \end{aligned} \quad (\text{A.1; A.3})$$

$$\begin{aligned} &\leq \frac{w_a^2}{4} + \frac{3w_b^2}{4} + \frac{(w_a + w_b)^2}{4} \\ &\leq \frac{w_a^2}{2} + w_b^2 + \frac{w_a w_b}{2} \leq (w_a + w_b)^2 = (f(a, c) + f(b, c))^2 \end{aligned} \quad (\text{A.4; A.5})$$

We now work on the last inequality. We want to show that $f(b, c) \leq f(a, b) + f(a, c)$. Consider the function $g_1(w_b) = (f(a, b) + f(a, c))^2 - f(b, c)^2$. We need to show that $g_1(w_b) \geq 0$ for w_b between w_a and $\sqrt{d_{\mathbb{X}}(a, b)^2 + w_a^2}$ due to relation (A.3).

$$\begin{aligned} g_1(w_b) &= \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} + w_a^2 + 2w_a \sqrt{\frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2}} - w_b^2 \\ &\geq \frac{3w_a^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} + 2w_a \sqrt{\frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2}} - \frac{d_{\mathbb{X}}(a, b)^2 + w_a^2}{2} \\ &\geq w_a^2 + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} + 2w_a \sqrt{\frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2}} - \frac{d_{\mathbb{X}}(a, b)^2}{4} = h_1(w_b) \end{aligned}$$

The function $h_1(w_b)$ is increasing with respect to w_b as $w_b \geq w_a$. Thus, we only need to show that $h(w_a) \geq 0$. Remark that (A.1) gives $d_{\mathbb{X}}(a, b)^2 \leq 2(d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2) = 2(w_a^2 + w_b^2)$. As $w_b = w_a$ here, we get $4w_a^2 \geq d_{\mathbb{X}}(a, b)^2$.

$$h_1(w_a) = w_a^2 + 2w_a \sqrt{w_a^2 + \frac{d_{\mathbb{X}}(a, b)^2}{4}} - \frac{d_{\mathbb{X}}(a, b)^2}{4} \geq w_a^2 + 2w_a^2 - w_a^2 \geq 0$$

Hence,

$$f(b, c) \leq f(a, c) + f(a, b)$$

Case 1.3: $w_b \geq w_c \geq w_a$

This case is very similar to the previous one. We now have the relations:

$$f(a, c) = w_c$$

$$f(b, c) = w_b$$

As in the previous section, we can deduce that:

$$f(a, c) \leq f(b, c) \leq f(b, c) + f(a, c)$$

Considering $f(a, b)$, the proof is almost the same replacing w_a by w_c :

$$\begin{aligned} f(a, b)^2 &= \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} \\ &\leq \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2}{4} + \frac{w_b^2 - w_a^2}{4} & (A.1; A.3) \\ &\leq \frac{w_a^2}{4} + \frac{3w_b^2}{4} + \frac{(w_c + w_b)^2}{4} & (A.4; A.5) \\ &\leq \frac{w_c^2}{2} + w_b^2 + \frac{w_c w_b}{2} & w_a \leq w_c \\ &\leq (w_c + w_b)^2 = (f(a, c) + f(b, c))^2 \end{aligned}$$

We consider $g_2(w_b) = (f(a, b) + f(a, c))^2 - f(b, c)^2$. We need to show that $g_2(w_b) \geq 0$ for w_b between w_c and $\sqrt{d_{\mathbb{X}}(a, b)^2 + w_a^2}$.

$$\begin{aligned} g_2(w_b) &= \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} + w_c^2 + 2w_c \sqrt{\frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2}} - w_b^2 \\ &\geq \frac{3w_a^2}{2} - \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} + 2w_a \sqrt{\frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2}} \\ &= g_1(w_b) \end{aligned}$$

Using the development on g_1 from Case 1.2, which do not use w_c , we have $g_2(w_b) \geq 0$ for

$$w_b \in [w_c, \sqrt{d_{\mathbb{X}}(a, b)^2 + w_a^2}] \supset [w_c, \sqrt{d_{\mathbb{X}}(a, b)^2 + w_a^2}]$$

Thus:

$$f(b, c) \leq f(a, c) + f(a, b).$$

One saturation We now release another pair. We hence have only one saturated pair. Without loss of generality, we assume that b is the common point to the unsaturated relations and that $w_c \geq w_a$. Then we have $f(a, c) = w_c$ and the following relations hold:

$$|w_b^2 - w_a^2| < d_{\mathbb{X}}(a, b)^2 \tag{A.6}$$

$$|w_c^2 - w_a^2| \geq d_{\mathbb{X}}(a, c)^2 \quad (\text{A.7})$$

$$|w_c^2 - w_b^2| < d_{\mathbb{X}}(b, c)^2 \quad (\text{A.8})$$

First, we consider $f(a, c) \leq f(a, b) + f(b, c)$ and introduce:

$$\begin{aligned} I &= (f(a, b) + f(b, c))^2 - f(a, c)^2 \\ &= w_b^2 + \frac{w_a^2 + w_c^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2 + d_{\mathbb{X}}(b, c)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} + \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} + 2f(a, b)f(b, c) - w_c^2 \\ &\geq w_b^2 - \frac{w_c^2}{2} + \frac{d_{\mathbb{X}}(b, c)^2}{4} + \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)} = g_3(w_c) \end{aligned}$$

In order to show that $I \geq 0$, we need to prove that $g_3(w_c) \geq 0$ for $w_c \in [0, \sqrt{d_{\mathbb{X}}(b, c)^2 + w_b^2}]$, using (A.8). g_3 is derivable and we obtain:

$$g_3'(w_c) = w_c \left(\frac{w_c^2 - w_b^2}{d_{\mathbb{X}}(b, c)} - 1 \right)$$

Due to relation (A.8) and $w_c \geq 0$, $g_3'(w_c) \leq 0$. The minimum for g_3 is thus reached for the greatest possible value of w_c and using (A.8) again, we have $w_c^2 \leq d_{\mathbb{X}}(b, c)^2 + w_b^2$.

$$g_3(d_{\mathbb{X}}(b, c)^2 + w_b^2) = w_b^2 - \frac{w_b^2 + d_{\mathbb{X}}(b, c)^2}{2} + \frac{d_{\mathbb{X}}(b, c)^2}{4} + \frac{(w_b^2 + d_{\mathbb{X}}(b, c)^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} = \frac{w_b^2}{2} \geq 0$$

Thus we can deduce that:

$$f(a, c) \leq f(a, b) + f(b, c)$$

For the second inequality, we consider the relation:

$$\begin{aligned} I &= (f(a, b) + f(a, c))^2 - f(b, c)^2 \\ &= \frac{w_a^2 + w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} + w_c^2 + 2w_c f(a, b) - \frac{w_b^2 + w_c^2}{2} - \frac{d_{\mathbb{X}}(b, c)^2}{4} - \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} \\ &\geq \frac{w_a^2}{2} + \frac{w_c^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} - \frac{d_{\mathbb{X}}(b, c)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} - \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} + w_c d_{\mathbb{X}}(a, b) = g_4 \end{aligned}$$

Consider g_4 as a function of $d_{\mathbb{X}}(b, c)$, we can derive:

$$g_4'(d_{\mathbb{X}}(b, c)) = \frac{(w_c^2 - w_b^2)^2}{2d_{\mathbb{X}}(b, c)^3} - \frac{d_{\mathbb{X}}(b, c)}{2}$$

Given the relation (A.8), we have $g_4'(d_{\mathbb{X}}(b, c)) \leq 0$ on our domain of interest $[\sqrt{|w_c^2 - w_b^2|}, d_{\mathbb{X}}(a, b) +$

$d_{\mathbb{X}}(a, c)$]. Thus:

$$\begin{aligned} g_4 &\geq g(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c)) \\ &= \frac{w_a^2}{2} + \frac{w_c^2}{2} + w_c d_{\mathbb{X}}(a, b) + \frac{d_{\mathbb{X}}(a, b)^2 - (d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} - \frac{(w_c^2 - w_b^2)^2}{4(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2} \\ &\geq \frac{w_a^2}{2} + \frac{w_c^2}{4} + \frac{w_c d_{\mathbb{X}}(a, b)}{2} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} - \frac{(w_c^2 - w_b^2)^2}{4(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2} = h_4 \quad (A.7) \end{aligned}$$

Now, considering h_4 as a function of w_a , h_4 is derivable and $h'_4(w_a) = w_a \left(1 - \frac{w_b^2 - w_a^2}{d_{\mathbb{X}}(a, b)^2}\right)$. Using (A.6), we have $h'_4(w_a) \geq 0$ on our domain. We thus consider w_a as small as possible. There are two cases to consider. If $w_b \leq d_{\mathbb{X}}(a, b)$ then the minimum is reached for $w_a = 0$. Otherwise, the minimum is reached for $w_a = \sqrt{w_b^2 - d_{\mathbb{X}}(a, b)^2}$.

Case 2.1: $w_b \leq d_{\mathbb{X}}(a, b)$

$$h_4(0) = \frac{w_c^2}{4} + \frac{w_c d_{\mathbb{X}}(a, b)}{2} + \frac{w_b^4}{4d_{\mathbb{X}}(a, b)^2} - \frac{(w_c^2 - w_b^2)^2}{4(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2} = i_4(w_b)$$

Deriving i_4 , we get:

$$i'_4(w_b) = w_b \left(\frac{w_b^2}{d_{\mathbb{X}}(a, b)^2} + \frac{w_c^2 - w_b^2}{(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2} \right) \geq 0$$

As $w_b \geq 0$, it suffices to check:

$$\begin{aligned} i_4(0) &= \frac{w_c^2}{4} + \frac{w_c d_{\mathbb{X}}(a, b)}{2} - \frac{w_c^4}{4(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2} \\ &\geq \frac{w_c^2}{4} \left[1 - \frac{w_c^2}{(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2} \right] \geq 0 \end{aligned}$$

Case 2.2: $w_b > d_{\mathbb{X}}(a, b)$

$$\begin{aligned} h_4(\sqrt{w_b^2 - d_{\mathbb{X}}(a, b)^2}) &\geq \frac{w_c^2}{4} + \frac{w_c d_{\mathbb{X}}(a, b)}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} - \frac{(w_c^2 - w_b^2)^2}{4(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2} \\ &\geq \frac{w_c^2}{4} + \frac{d_{\mathbb{X}}(a, c) d_{\mathbb{X}}(a, b)}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} - \frac{(d_{\mathbb{X}}(a, b) + d_{\mathbb{X}}(a, c))^2}{4} \quad (A.7; A.8; A.1) \\ &\geq 0 \quad (A.7) \end{aligned}$$

We thus have proved that:

$$f(b, c) \leq f(a, b) + f(a, c)$$

Let us now get to the last inequality. The proof is similar to the previous one:

$$\begin{aligned}
I &= (f(b, c) + f(a, c))^2 - f(a, b)^2 \\
&= w_c^2 + 2w_c f(b, c) + \frac{w_c^2 + w_b^2}{2} + \frac{d_{\mathbb{X}}(b, c)^2}{4} + \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} - \frac{w_a^2 + w_b^2}{2} - \frac{d_{\mathbb{X}}(a, b)^2}{4} - \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} \\
&\geq w_c^2 + d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c) + \frac{d_{\mathbb{X}}(b, c)^2}{4} + \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} - \frac{d_{\mathbb{X}}(a, b)^2}{4} - \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2} = g_5
\end{aligned}$$

g_5 considered as a function of $d_{\mathbb{X}}(a, b)$ is derivable and:

$$g_5'(d_{\mathbb{X}}(a, b)) = \frac{d_{\mathbb{X}}(a, b)}{2} \left(\frac{(w_b^2 - w_a^2)^2}{d_{\mathbb{X}}(a, b)^4} - 1 \right) \leq 0$$

using the relation (A.6). Thus:

$$\begin{aligned}
g_5(d_{\mathbb{X}}(a, b)) &\geq g_5(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c)) \\
&= w_c^2 + d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c) + \frac{d_{\mathbb{X}}(b, c)^2 - (d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2}{4} + \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} - \frac{(w_b^2 - w_a^2)^2}{4(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2} \\
&\geq \frac{3w_c^2}{4} + \frac{d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c)}{2} + \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} - \frac{(w_b^2 - w_a^2)^2}{4(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2} = h_5 \quad (A.7)
\end{aligned}$$

Consider h_5 as a function of w_c .

$$h_5'(w_c) = w_c \left(\frac{3}{2} + \frac{w_c^2 - w_b^2}{d_{\mathbb{X}}(b, c)^2} \right)$$

Knowing that $w_c^2 \geq w_b^2 - d_{\mathbb{X}}(b, c)^2$ due to (A.8) and that $w_c \geq 0$, we have $h'(w_c) \geq 0$. We thus have two cases to consider:

Case 3.1: $w_b^2 - d_{\mathbb{X}}(b, c)^2 \geq 0$ We look at $h_5(\sqrt{w_b^2 - d_{\mathbb{X}}(b, c)^2})$:

$$\begin{aligned}
h_5(\sqrt{w_b^2 - d_{\mathbb{X}}(b, c)^2}) &\geq \frac{w_c^2}{2} + \frac{w_b^2}{4} - \frac{d_{\mathbb{X}}(b, c)^2}{4} + \frac{d_{\mathbb{X}}(b, c)^2}{4} - \frac{|w_b^2 - w_a^2|}{4} \\
&\geq \frac{w_c^2}{2} + \frac{w_b^2}{4} - \frac{w_b^2}{4} - \frac{w_a^2}{4} \geq \frac{w_c^2}{4} \geq 0
\end{aligned}$$

Case 3.2: $w_b^2 - d_{\mathbb{X}}(b, c)^2 \leq 0$ We know that $w_c^2 \geq d_{\mathbb{X}}(a, c)^2 + w_a^2$ due to (A.7). Thus:

$$\begin{aligned}
 h_5(w_c) &\geq h_5(d_{\mathbb{X}}(a, c)^2 + w_a^2) \\
 &= \frac{3d_{\mathbb{X}}(a, c)^2}{4} + \frac{3w_a^2}{4} + \frac{d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c)}{2} + \frac{(d_{\mathbb{X}}(a, c)^2 + w_a^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} - \frac{(w_b^2 - w_a^2)^2}{4(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2} \\
 &\geq \frac{3d_{\mathbb{X}}(a, c)^2}{4} + \frac{3w_a^2}{4} + \frac{d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c)}{2} + \frac{(d_{\mathbb{X}}(a, c)^2 + w_a^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} - \frac{|w_b^2 - w_a^2|}{4} \\
 &\geq \frac{3d_{\mathbb{X}}(a, c)^2}{4} + \frac{w_a^2}{2} + \frac{d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c)}{2} + \frac{(d_{\mathbb{X}}(a, c)^2 + w_a^2 - w_b^2)^2}{4d_{\mathbb{X}}(b, c)^2} - \frac{w_b^2}{4}
 \end{aligned} \tag{A.6}$$

If $w_b^2 \leq w_c^2 = d_{\mathbb{X}}(a, c)^2 + w_a^2$, then $h_5(w_c) \geq 0$. Otherwise we can use the relation (A.8) to state that $w_c^2 - w_b^2 > -d_{\mathbb{X}}(b, c)^2$. Thus:

$$\begin{aligned}
 h(w_c) &\geq \frac{3d_{\mathbb{X}}(a, c)^2}{4} + \frac{w_a^2}{2} + \frac{d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c)}{2} - \frac{d_{\mathbb{X}}(a, c)^2 + w_a^2 - w_b^2}{4} - \frac{w_b^2}{4} \\
 &\geq \frac{d_{\mathbb{X}}(a, c)^2}{2} + \frac{w_a^2}{4} \geq 0
 \end{aligned}$$

This conclude our proof and we have:

$$f(a, b) \leq f(b, c) + f(a, c)$$

No saturation Finally, we need to consider the case where no pair is saturated. This means that the three following relations are verified:

$$|w_b^2 - w_a^2| < d_{\mathbb{X}}(a, b)^2 \tag{A.9}$$

$$|w_c^2 - w_a^2| < d_{\mathbb{X}}(a, c)^2 \tag{A.10}$$

$$|w_c^2 - w_b^2| < d_{\mathbb{X}}(b, c)^2 \tag{A.11}$$

Remark that the three points are distinct. We want to show that $f(a, b) \leq f(b, c) + f(a, c)$. Let us consider the square of the relation:

$$f(a, b)^2 = \frac{w_a^2}{2} + \frac{w_b^2}{2} + \frac{d_{\mathbb{X}}(a, b)^2}{4} + \frac{(w_b^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, b)^2}$$

and:

$$(f(b, c) + f(a, c))^2 = f(b, c)^2 + f(a, c)^2 + 2f(a, c)f(b, c)$$

Remark that $f(a, c)f(b, c) \geq \frac{d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c)}{4}$. Hence,

$$(f(b, c) + f(a, c))^2 \geq \frac{w_b^2 + w_c^2}{2} + \frac{d_{\mathbb{X}}(b, c)^2}{4} + \frac{(w_c^2 - w_b^2)^2}{4d_{\mathbb{X}}(a, b)^2} + \frac{w_a^2 + w_c^2}{2} + \frac{d_{\mathbb{X}}(a, c)^2}{4} + \frac{(w_c^2 - w_a^2)^2}{4d_{\mathbb{X}}(a, c)^2} + 2 \frac{d_{\mathbb{X}}(b, c)}{2} \frac{d_{\mathbb{X}}(a, c)}{2}$$

Multiplying the relation by 4 and using $w_c^2 \geq 0$, it is then sufficient to prove that:

$$d_{\mathbb{X}}(a, b)^2 + \frac{(w_b^2 - w_a^2)}{d_{\mathbb{X}}(a, b)^2} \leq d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2 + 2d_{\mathbb{X}}(a, c)d_{\mathbb{X}}(b, c) + \frac{(w_c^2 - w_b^2)^2}{d_{\mathbb{X}}(b, c)^2} + \frac{(w_c^2 - w_a^2)^2}{d_{\mathbb{X}}(a, c)^2}$$

We fix:

$$g_6(d_{\mathbb{X}}(a, b)) = (d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2 - d_{\mathbb{X}}(a, b)^2 + \frac{(w_c^2 - w_b^2)^2}{d_{\mathbb{X}}(b, c)^2} + \frac{(w_c^2 - w_a^2)^2}{d_{\mathbb{X}}(a, c)^2} - \frac{(w_b^2 - w_a^2)^2}{d_{\mathbb{X}}(a, b)^2}$$

Given relations (A.9) and (A.1), we need to study g_6 for $d_{\mathbb{X}}(a, b)$ in the interval

$I =]\sqrt{|w_b^2 - w_a^2|}, d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c)]$ and g_6 is derivable. Thus:

$$g'_6(d_{\mathbb{X}}(a, b)) = 2 \frac{(w_b^2 - w_a^2)^2}{d_{\mathbb{X}}(a, b)^3} - 2d_{\mathbb{X}}(a, b)$$

Remark that g'_6 is equal to 0 for $d_{\mathbb{X}}(a, b) = \sqrt{|w_b^2 - w_a^2|}$ and is negative on I . Thus the minimum of g_6 is reached for $d_{\mathbb{X}}(a, b) = d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c)$.

$$\begin{aligned} g_6(d_{\mathbb{X}}(a, b)) &\geq g_6(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c)) \\ &= \frac{(w_c^2 - w_b^2)^2}{d_{\mathbb{X}}(b, c)^2} + \frac{(w_c^2 - w_a^2)^2}{d_{\mathbb{X}}(a, c)^2} - \frac{(w_b^2 - w_a^2)^2}{(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2} = h_6(w_c) \end{aligned}$$

Let us study h_6 on \mathbb{R}_+ . h_6 is differentiable and:

$$h'_6(w_c) = 2w_c \left(\frac{w_c^2 - w_b^2}{d_{\mathbb{X}}(b, c)^2} + \frac{w_c^2 - w_a^2}{d_{\mathbb{X}}(a, c)^2} \right)$$

$h'_6 = 0$ if $w_c = 0$ or $w_c^2 = \frac{w_a^2 d_{\mathbb{X}}(b, c)^2 + w_b^2 d_{\mathbb{X}}(a, c)^2}{d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2} = x^2$ and h'_6 is negative on $[0, x]$ and positive for $w_c > x$. Thus the minimum of h is reached for $w_c = x$.

$$\begin{aligned} h_6(x) &= \frac{\left(\frac{w_a^2 d_{\mathbb{X}}(b, c)^2 - w_b^2 d_{\mathbb{X}}(b, c)^2}{d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2} \right)^2}{d_{\mathbb{X}}(b, c)^2} + \frac{\left(\frac{w_b^2 d_{\mathbb{X}}(a, c)^2 - w_a^2 d_{\mathbb{X}}(a, c)^2}{d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2} \right)^2}{d_{\mathbb{X}}(a, c)^2} - \frac{(w_b^2 - w_a^2)^2}{(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c))^2} \\ &= (w_b^2 - w_a^2)^2 \left(\frac{d_{\mathbb{X}}(b, c)^2}{(d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2)^2} + \frac{d_{\mathbb{X}}(a, c)^2}{(d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2)^2} - \frac{1}{(d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2)} \right) \\ &= (w_b^2 - w_a^2)^2 \left(\frac{1}{d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2} - \frac{1}{d_{\mathbb{X}}(a, c)^2 + d_{\mathbb{X}}(b, c)^2} \right) = 0 \end{aligned}$$

Consequently, $h_6(w_c) \geq 0$ for all w_c in \mathbb{R}_+ and this means $g_6(d_{\mathbb{X}}(a, c) + d_{\mathbb{X}}(b, c)) \geq 0$. We can hence conclude that:

$$f(a, b) \leq f(a, c) + f(b, c)$$

■

B Résumé en français

B.1 Introduction

La collecte de données est aujourd'hui une partie intégrante de la vie quotidienne. Des entreprises qui collectent des informations sur leurs clients et leurs employés aux services de renseignement, en passant par les instituts de sondage, la quête d'information est partout présente. Les données sont traitées comme une ressource qu'il faut amasser et utiliser afin d'améliorer son efficacité. Les gouvernements espèrent encourager la croissance économique par ce biais, comme en témoignent les nombreuses initiatives de type données ouvertes [1, 2, 3]. Cependant, les données elles-mêmes sont inutiles. Il est nécessaire de les interpréter et de les transformer en information.

L'interprétation est souvent réalisée grâce à une visualisation. Dans le cas d'objets en deux ou trois dimensions, l'être humain est capable de les visualiser. Cependant, les données sont souvent obtenues comme un nuage de points dans un espace de grande dimension. Par exemple, une image en niveaux de gris est un point dans un espace dont la dimension est son nombre de pixels. De tels espaces de grande dimension sont impossibles à visualiser directement et nécessitent un prétraitement avant leur interprétation.

De nombreux problèmes font partie de cette interprétation. Le clustering [41, 59, 78, 95] et la segmentation [98] tentent de séparer les points en différents groupes. La reconstruction essaye de construire un objet continu à partir des points de données, généralement sous la forme d'une triangulation [21, 47]. La réduction de dimension projette les données sur un espace de plus petite dimension en utilisant les paramètres les plus pertinents pour décrire ces données [62, 93], ce qui peut rendre la visualisation et l'analyse plus aisées.

Dans cette thèse, nous considérons l'analyse topologique des données et plus précisément l'inférence de topologie. Nous cherchons à extraire de la structure à partir des données en inférant la topologie sous-jacente. Cette information peut nous guider dans la résolution des problèmes mentionnés ci-dessus. Si nous connaissons le nombre de composantes connexes alors nous connaissons le nombre de groupes que nous devrions obtenir par clustering. Si nous connaissons la dimension intrinsèque des données alors nous connaissons la dimension de l'espace dont nous avons besoin pour une réduction de dimension. Ces dernières années, l'homologie persistante a été un des outils les plus populaires de l'analyse topologique des données. Cet outil analyse les données à toutes les échelles.

Il ne faut pas oublier que les données disponibles en pratique sont presque toujours bruitées,

que ce soit à cause d'erreurs de mesures ou de modélisations imparfaites. L'analyse topologique des données a besoin d'être robuste au bruit. Les algorithmes existants fonctionnent généralement bien lorsque le bruit est borné. Cependant, la présence de valeurs aberrantes est courante dans les données. Un appareil de mesure défectueux ou une erreur peuvent créer des points qui n'ont aucune relation avec le reste des données et qui sont difficiles à gérer, entraînant la malfonction des algorithmes.

Récemment, l'homologie persistante a été utilisée dans de nombreux domaines. Une première application est la définition de signatures pour les données. Dans ce cadre, l'homologie persistante fournit une information topologique qui discrimine différentes classes de phénomènes. Cela a été utilisé pour classifier des images de différentes pathologies [4, 38], pour analyser des électroencéphalogrammes [97] ou différencier des formes tridimensionnelles [23]. De plus, elle peut fournir une méthode pour la segmentation et le clustering des données [33, 85, 90]. L'étape suivante est la recherche de motifs pour la détection et l'identification de phénomènes, ce qui a été appliqué à des images [77], à des sous-types de cancer [84] et à des motifs cycliques du génome [46].

L'homologie persistante fournit également une manière de mieux comprendre la structure des objets et de la visualiser, de la structure de la matière en astrophysique [91, 92] aux matériaux à granularité dense [75], réseaux complexes [72, 86] et systèmes dynamiques [10]. En biologie, elle peut expliquer la compressibilité des protéines [65] et décrire les structures de racines [57]. Elle a également été utilisée pour étudier la propagation de gènes codant des résistances antibiotiques [58]. La reconstruction elle-même peut être réalisée [35], fournissant une structure pour le pistage [9] et une visualisation des structures corticales [76, 89].

Le produit de l'homologie persistante est généralement un diagramme de persistance, structure qui est peu adaptée à un contexte statistique. Par exemple, nous ne savons pas définir la moyenne de deux diagrammes. Cependant, l'introduction de paysages persistants [16] a rendu possible l'utilisation d'analyses statistiques pour la topologie, comme cela a été fait pour des données orthodontiques [64, 70].

B.2 Homologie persistante

Étant donné une famille d'espaces topologiques indexés par un paramètre $\alpha \in \mathbb{R}$, $\mathcal{F} = \{F_\alpha\}$, l'homologie persistante étudie l'évolution de la topologie de ces espaces tandis que α évolue de $-\infty$ à ∞ . En analyse de données, la méthode la plus usuelle pour construire une suite d'espaces topologiques est de faire grossir des boules. Étant donné un nuage de point P , cela signifie considérer les sous-niveaux de la distance à P . Supposant que P est un échantillonnage d'un objet sous-jacent K , nous espérons que certains des sous-niveaux ont la même topologie que K . La topologie des sous-niveaux est souvent stable pour un intervalle de valeurs. Dans cette thèse, nous ne considérons que des paramètres définis sur \mathbb{R} .

Par topologie, nous entendons homologie. Intuitivement, cela correspond aux composantes connexes en dimension 0, aux trous ou cycles en dimension 1, aux cavités en dimension 2 et ainsi de suite. Considérons la figure B.1. Nous avons un ensemble de points P qui échantillonne avec bruit les arêtes S d'un carré. Nous voulons inférer la topologie de S qui a une composante connexe et un cycle. Pour $\alpha < 0$, le sous-niveau de la distance à P est vide.

Quand $\alpha = 0$, le sous-niveau est exactement P et nous avons donc 14 composantes connexes, une pour chaque point de P . Quand α croît et que les boules grossissent, nous obtenons un sous-niveau qui a la même topologie que S . Remarquons que cette topologie est stable pour un certain intervalle de valeurs de α .

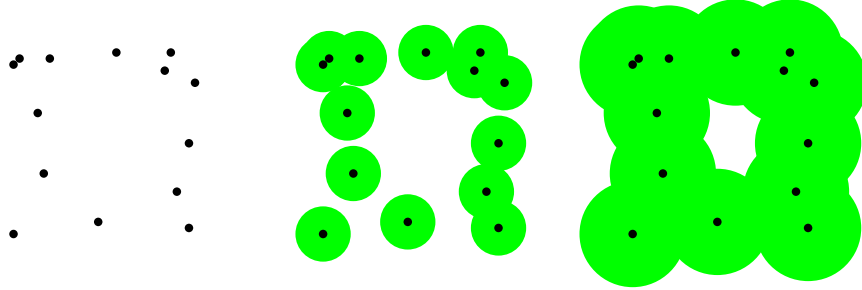


Figure B.1 – Croissance des boules pour la persistance

L'information topologique obtenue en utilisant la persistance est usuellement représentée par un diagramme de persistance. Un élément topologique, par exemple un cycle, apparaît dans un des espaces topologiques $F_\alpha \in \mathcal{F}$. α est appelée la *date de naissance* de l'élément topologique. Cet élément existe alors dans un certain nombre de F_γ tels que $\alpha < \gamma < \beta$ et disparaît dans les F_δ pour $\delta > \beta$ pour un certain β . β est appelé la *date de mort* de l'élément topologique. Notons qu'un des éléments de dimension 0, id est, une composante connexe ne meurt pas et a ainsi une date de mort infinie. Le diagramme de persistance de dimension d de \mathcal{F} est le multi-ensemble composé des paires de points (x, y) où x est la date de naissance d'un élément topologique de dimension d et y est la date de mort de cet élément. Les diagrammes de persistance peuvent être représentés, soit par un multi-ensemble de R^2 , soit par un code-barre. Dans le premier cas, chaque paire (x, y) est représentée par un point. Dans le second cas, (x, y) est représentée par une barre commençant en x et finissant en y . La figure B.2 montre les deux représentations obtenues pour les éléments topologiques de dimension 1, les cycles, dans l'exemple de la figure B.1.

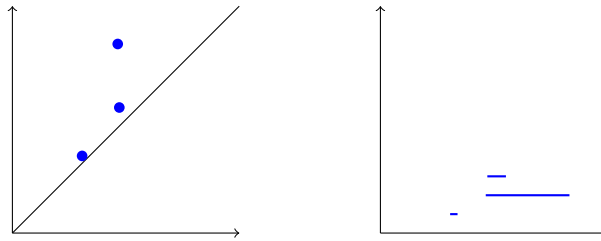


Figure B.2 – Diagramme de persistance en dimension 1

L'idée derrière la persistance est que les éléments topologiques correspondant à la structure de l'objet sous-jacent K sont stables sur un intervalle de valeurs du paramètres. Ainsi, ils ont

une plus grande durée de vie que les éléments topologiques créés par le bruit. Dans notre exemple, nous voyons qu'un élément de dimension 1 a une plus grande durée de vie que les autres. Il correspond au cycle de S . Les deux barres plus petites représentent de petits cycles, assimilés à du bruit, apparaissant lors de la croissance des boules.

La persistance est multi-échelle. Nous considérons l'ensemble des valeurs de α et nous observons ainsi les données à toutes les échelles. Cela signifie que nous pouvons détecter et inférer la topologie d'objets qui ont une topologie qui diffère selon l'échelle considérée. Par exemple, le nuage de point de la figure B.3 échantillonne une spirale enroulée sur un tore. Quand nous regardons ce nuage de très près, nous avons un objet de dimension 1, la spirale. À une distance intermédiaire, nous avons le tore qui est un objet de dimension 2. L'homologie persistante est capable d'analyser correctement cette différence de topologie dépendant de l'échelle.

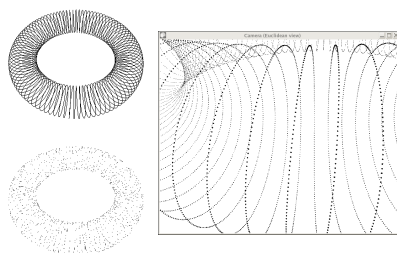


Figure B.3 – Spirale enroulée sur un tore

Un bon diagramme de persistance pour l'inférence topologique est un diagramme où le rapport entre la plus petite durée de vie d'un élément topologique pertinent et la plus longue durée de vie d'un élément dû au bruit, appelé *fossé*, est grand. Les diagrammes de persistance sont stables par rapport à de petites variations de la fonction utilisée pour définir les sous-niveau. Lorsque la distance à P approxime la distance à K , nous obtenons un bon diagramme. Le calcul des diagrammes de persistance n'échappe pas à la *malédiction de la dimension*. Le calcul fonctionne bien en petite dimension mais pas en grande dimension en raison d'une explosion de la complexité. La technique classique pour calculer les diagrammes de persistance d'une famille \mathcal{F} d'unions de boules est de construire une famille croissante de complexes simpliciaux \mathcal{G} . Un complexe simplicial est un ensemble de points, arêtes, triangles, tétraèdres et ainsi de suite. La famille \mathcal{G} approxime la topologie de \mathcal{F} . L'algorithme classique [57] calcule l'homologie persistante en temps $O(N^3)$, où N est le nombre de simplexes du complexe maximum dans \mathcal{G} . Cependant, si \mathcal{F} est décrit en utilisant n boules dans un espace de dimension d , nous devons construire le complexe simplicial maximum de dimension d . Sa taille est $\binom{n}{d}$. Ainsi, la complexité sera de l'ordre de $O(n^{3d})$ ce qui rend impossible son utilisation en grande dimension.

Des approches récentes cherchent à obtenir une complexité qui dépend de la dimension intrinsèque de l'objet au lieu de la dimension extrinsèque. Par exemple, cela a été obtenu pour les complexes de Vietoris-Rips [88]. Cela signifie qu'un objet de petite dimension plongé dans un espace de dimension plus grande peut être analysé sans payer le coût de complexité de l'espace ambiant. C'est une façon de contourner la malédiction de la dimension.

B.3 Le problème du bruit aberrant

La présence de valeurs aberrantes crée des problèmes lors du calcul des diagrammes de persistance. Considérons le 1-squelette d'un cube, c'est-à-dire l'ensemble de ses arêtes. En entrée, nous avons un ensemble de point qui échantillonne le squelette et contient quatre points aberrants situés au centre de quatre des faces du cubes, de telle sorte que les deux faces vides soient opposées, comme dans la figure B.4. Ces points de bruit perturbent le diagramme de persistance de manière significative et rendent l'inférence topologique impossible.

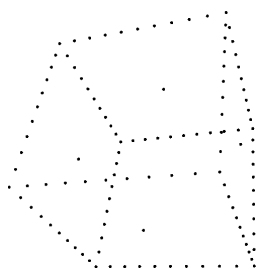


Figure B.4 – Échantillonnage du squelette d'un cube avec points aberrants

Nous souhaitons inférer le diagramme de persistance du squelette de cube de la figure B.5, c'est-à-dire le diagramme de persistance des sous-niveaux de la distance au squelette du cube. Cet objet a une unique composante connexe naissant en 0 et existant pour toute les valeurs positives du paramètre α . Au début, nous avons 5 éléments topologiques de dimension 1, ou cycles, car le cube a six faces et l'une d'entre elle est la somme des cinq autres. Lorsque le paramètre grandit, les faces sont remplies et les éléments de dimension 1 disparaissent, remplacés par un élément de dimension 2 correspondant à la cavité à l'intérieur du cube.

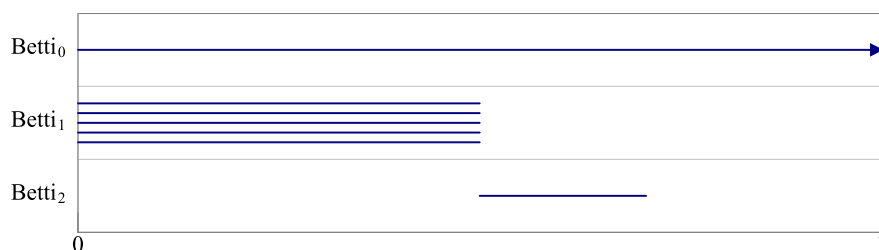


Figure B.5 – Diagramme de persistance du squelette de cube

La présence de bruit aberrant brouille le diagramme de persistance. Le diagramme de persistance obtenu pour la filtration des sous-niveaux de la distance au nuage de points est donné dans la figure B.6. Notons que le diagramme de dimension 1 a maintenant un plus petit fossé mais nous pouvons encore inférer la bonne structure. Cependant, en dimension 2, l'homologie persistante est complètement différente et le fossé est de 1, ce qui signifie que nous ne pouvons plus différencier le signal du bruit. Nous avons deux éléments topologiques de même durée de vie. Chacun correspond à la moitié du cube. Quand les faces sont remplies par la croissance des boules, une connexion s'opère également au milieu du cube, créée par

les quatre points aberrants.

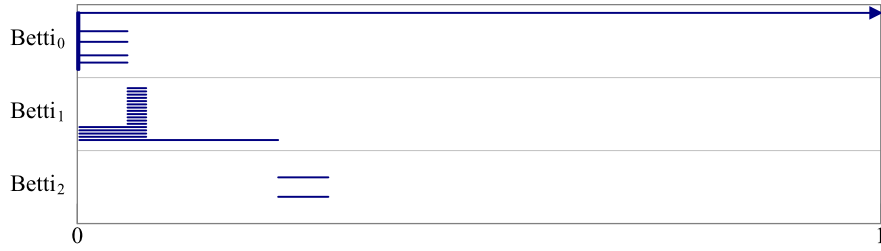


Figure B.6 – Diagramme de persistance obtenu en utilisant l'échantillonnage avec bruit aberrant

Inférer le diagramme complètement est particulièrement utile lorsque l'homologie persistante est utilisée pour obtenir des signatures. Des différences dans les parties tardives des diagrammes peuvent fournir des informations intéressantes pour discriminer les objets. Cependant, les points aberrants peuvent complètement changer l'aspect des diagrammes.

B.4 La distance à la mesure

Afin de gérer le bruit et en particulier les points aberrants, l'idée est de remplacer la distance au nuage de points P par une autre fonction. Une telle fonction doit avoir deux propriétés. Elle doit être stable par rapport à de petites variations dans les données et ses sous-niveaux doivent être facilement calculables.

Nous utilisons la distance à la mesure. Étant donné un ensemble de n points P dans un espace métrique \mathbb{X} et un paramètre de masse $m = \frac{k}{n}$, où k est un entier, la distance à la mesure empirique μ sur P est la fonction définie sur \mathbb{X} par

$$d_{\mu,m}(x) = \sqrt{\frac{1}{k} \sum_{i=1}^k d_{\mathbb{X}}(x, p_i(x))^2}$$

où $p_i(x)$ est le $i^{\text{ème}}$ plus proche voisin de x dans P et $d_{\mathbb{X}}(x, p_i(x))$ est la distance entre x et $p_i(x)$. Dans un contexte plus général, la distance à une mesure μ en x quantifie le coût du meilleur plan de transport pour amener une masse m depuis μ jusqu'à x . Cette fonction est stable [25] et peut-être aisément calculée dans les espaces euclidiens [67]. C'est donc un bon candidat pour inférer la topologie.

Revenons au squelette de cube. Le diagramme de persistance de la distance à la mesure empirique sur P pour une masse correspondant à 5 points est donnée dans le figure B.7. Le diagramme contient encore du bruit mais il y a une nette différence entre les durées de vie des éléments pertinents et des éléments dus au bruit en dimension 1 et 2.

L'utilisation du diagramme de persistance de la distance à la mesure sur des données réelles est confrontée à un problème majeur. En dehors des espaces euclidiens, les sous-niveaux de la distance à la mesure ne sont pas nécessairement calculables. Dans les espaces euclidiens, les sous-niveaux sont des unions de boules [67]. Cependant, le nombre de boules nécessaires

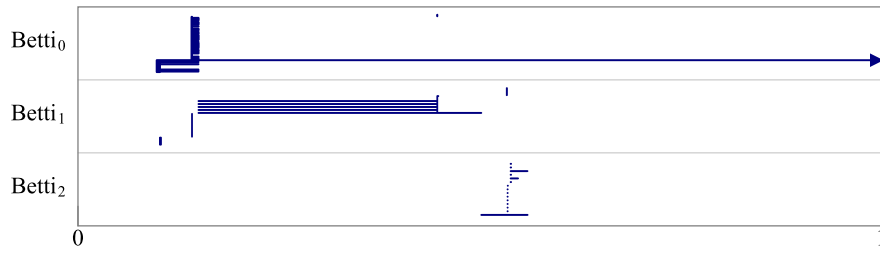


Figure B.7 – Diagramme obtenu en utilisant la distance à la mesure empirique pour le squelette de cube

pour décrire cette union est identique au nombre de cellules de Voronoi d'ordre k qui sont non-vides et qui peut être de l'ordre de $O\left(n^{\lfloor \frac{d+1}{2} \rfloor} k^{\lceil \frac{d+1}{2} \rceil}\right)$ [39]. Il est donc nécessaire d'approximer les sous-niveaux avant de pouvoir les utiliser en pratique. Une approximation avec un nombre linéaire de boules a été proposée dans [67] et des bornes inférieures sur le nombre minimal de boules nécessaires sont données dans [80]. Malheureusement, ces résultats sont limités aux espaces euclidiens et ne peuvent être étendus à d'autres espaces métriques.

Cette approximation n'est pas suffisante pour traiter des données en grande dimension. La grande taille des complexes simpliciaux utilisés pour le calcul des diagrammes de persistance existe toujours. Les résultats de [88] supposent que pour un $F_\alpha \in \mathcal{F}$, toutes les boules ont même rayon. Les boules dans les sous-niveaux de $d_{\mu,m}$ ne vérifient pas cette propriété et la méthode de [88] a besoin d'être adaptée.

B.5 Contenu de la thèse

Cette thèse étudie la complexité du calcul de l'homologie persistante et la manière dont nous pouvons gérer le bruit et en particulier les valeurs aberrantes. Notre but est de rendre l'homologie persistante robuste au bruit et utilisable en pratique pour des données de petite dimension intrinsèque, éventuellement plongées dans un espace de grande dimension. Nous proposons une méthode d'approximation des diagrammes de persistance de la distance à la mesure. Nous introduisons également de nouvelles conditions d'échantillonnage mieux adaptées à l'utilisation de la distance à la mesure. Cela nous permet d'élargir l'ensemble des applications possibles.

B.5.1 Distance à la mesure

La distance à une mesure μ était définie dans les espaces euclidiens [25]. Sa stabilité est garantie si deux mesures sont proches en distance de Wasserstein. La définition et la stabilité de la distance à la mesure peuvent être étendues trivialement au cas des espaces métriques.

Théorème B.1 ¹ Soient μ et ν deux mesures de probabilité sur un espace métrique \mathbb{X} et $m \in]0, 1[$ alors :

$$\|d_{\mu,m} - d_{\nu,m}\|_\infty \leq \frac{1}{\sqrt{m}} W_2(\mu, \nu).$$

¹Théorème 3.15

Une fois établie la stabilité de la distance à une mesure, la question de l'identifiabilité se pose naturellement. Cette question est notamment cruciale pour l'utilisation de la distance à la mesure pour la construction de signatures. Une bonne signature doit être capable de discriminer entre différents objets. Nous fournissons de nouveaux résultats dans les espaces euclidiens dans la section 3.5. En particulier, la distance à la mesure considérée pour toutes les valeurs de m permet d'identifier une mesure de probabilité.

Théorème B.2² Soient μ et ν deux mesures de probabilités sur \mathbb{R}^d , alors :

$$\left(\forall x \in \mathbb{R}^d, \forall m \in [0, 1[, d_{\mu, m}(x) = d_{\nu, m}(x) \right) \Leftrightarrow \mu = \nu.$$

La restriction à des mesures empiriques définies sur un nuage fini de points permet de réduire les condition sur m .

Théorème B.3³ Soient $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction, $d \geq 2$, et $m \in [0, \frac{1}{2}[$. S'il existe un nuage fini de points P tel que la mesure empirique μ sur P vérifie $f = d_{\mu, m}$ et pour tout $x \in \mathbb{R}^d$, $\mu(x) \leq m$, alors la mesure μ est définie de manière unique.

L'utilisation de la distance à la mesure dans le cadre de l'homologie persistante nécessite également de pouvoir facilement calculer ses sous-niveaux. Lorsque la mesure considérée est une mesure empirique sur un nuage de n points, des travaux précédents ont proposé une approximation par une distance de puissance, appelée *witnessed k -distance* [67]. Cette approximation a une taille linéaire par rapport au nombre de points du nuage original, c'est-à-dire que les sous-niveaux sont décrits par une union d'au plus n boules. Cependant, elle repose sur l'existence de barycentres, ce qui limite son utilisation aux espaces euclidiens. De plus elle introduit de nouveaux points comme support de l'approximation ce qui peut s'avérer problématique pour l'analyse de champs scalaires. Enfin, elle n'est pas applicable au cas des mesures dont le support n'est pas fini. Reprenant l'idée d'une approximation par une distance de puissance, nous introduisons une nouvelle fonction d'approximation.

Définition B.4⁴ Soient μ une mesure de probabilité sur un espace métrique \mathbb{X} et $m \in [0, 1[$ un paramètre de masse. Étant donné un sous-ensemble P de \mathbb{X} , nous définissons $d_{\mu, m}^P$ Comme la distance de puissance associé à $(P, d_{\mu, m})$:

$$d_{\mu, m}^P(x) = \sqrt{\min_{p \in P} d_{\mu, m}(p)^2 + d_{\mathbb{X}}(p, x)^2}.$$

Cette fonction est définie dans tout espace métrique. Nous montrons que si l'ensemble P est un échantillonnage suffisamment dense du support de la mesure μ alors $d_{\mu, m}^P$ approxime $d_{\mu, m}$. En particulier, si μ est une mesure empirique sur un nuage de n points et que P est ce nuage alors nous obtenons une fonction dont les sous-niveaux sont déniés par une union d'au plus n boules, comme la *witnessed k -distance*. Cette fonction approxime la distance à

²Théorème 3.17

³Théorème 3.24

⁴Définition 4.11

μ dans n'importe quel espace métrique. De plus, les techniques de preuve utilisées permettent d'améliorer les bornes d'approximation de la witnessed k -distance $d_{\mu,m}^W$ qui deviennent équivalentes à celles de $d_{\mu,m}^P$.

Théorème B.5 ⁵ Soient P un ensemble fini de points d'un espace métrique \mathbb{X} , μ la mesure empirique sur P et $m \in [0, 1[$ un paramètre de masse alors les bornes suivantes sont optimales :

$$\frac{1}{\sqrt{2}} d_{\mu,m} \leq d_{\mu,m}^P \leq \sqrt{5} d_{\mu,m}.$$

De plus, si l'espace \mathbb{X} est euclidien alors les bornes deviennent :

$$\frac{1}{\sqrt{2}} d_{\mu,m} \leq d_{\mu,m}^P \leq \sqrt{3} d_{\mu,m}$$

et pour la witnessed k -distance :

$$d_{\mu,m} \leq d_{\mu,m}^W \leq \sqrt{6} d_{\mu,m}.$$

B.5.2 Structures de données pour la persistance des distances de puissance

L'approximation de la distance à la mesure a permis de se ramener à une distance de puissance définie par un nombre linéaire de boules. Cependant, la difficulté du calcul de son diagramme de persistance existe toujours. Une adaptation naturelle de la méthode basée sur le nerf des boules, permet de définir une version pondérée de la filtration de Vietoris-Rips R_α . Cette structure est stable et induit une métrique sur le complexe ainsi défini.

Théorème B.6 ⁶ Soient (P, w) un nuage de points pondéré et f la fonction induite par le Rips pondéré. Alors,

$$\forall a, b, c \in P, f(a, c) \leq f(a, b) + f(b, c)$$

Cependant, la taille de la filtration conserve les défauts du cas non pondéré. En particulier, elle explose lorsque la dimension augmente. Nous adaptons la méthode de Rips parcimonieux S_α développée par Don Sheehy [88]. Remplacer la métrique euclidienne par la nouvelle métrique du Rips pondéré ne permet pas de conserver les garanties de taille que nous recherchons. Il est nécessaire de construire le Rips parcimonieux dans le cas non pondéré puis d'intersecter ce dernier avec le Rips pondéré R_α . Cela définit un complexe T_α pour tout paramètre α et une filtration $\{T_\alpha\}$. Le calcul de ce complexe est relativement aisé. De plus, une fois calculé la filtration $\{S_\alpha\}$, un changement des poids définissant R_α correspond à un simple réordonnancement de l'apparition des arêtes dans $\{S_\alpha\}$. En définitive, nous obtenons une filtration qui conserve les propriétés de parcimonie de $\{S_\alpha\}$, à savoir une taille en $O(C^l n)$ où n est la taille de P , C une constante et l la dimension intrinsèque de l'objet échantillonné par P . De plus, le diagramme de persistance de $\{T_\alpha\}$ approxime bien celui de $\{R_\alpha\}$:

⁵Corollaire 4.13 et théorèmes 4.7 et 4.15

⁶Théorème 5.12

Théorème B.7 ⁷ Soit (P, w) , un sous-ensemble pondéré fini d'un espace métrique \mathbb{X} dont les poids sont t -Lipschitz. Soit $\epsilon \in]0, 1[$ un paramètre fixé pour la construction de la filtration parcimonieuse $\{T_\alpha\}$. Alors,

$$d_B^{\log}(\text{Dgm}(\{T_\alpha\}), \text{Dgm}(\{R_\alpha\})) \leq \ln \left(\frac{1 + \sqrt{1 + t^2 \epsilon}}{1 - \epsilon} \right).$$

La combinaison de l'approximation de la distance à la mesure et l'utilisation de la filtration du Rips pondéré parcimonieux permet de rendre le calcul du diagramme de persistance utilisable en pratique pour des données de petite dimension intrinsèque mais éventuellement plongées dans un espace de grande dimension.

B.5.3 Nouvelles conditions d'échantillonnage pour la distance à la mesure

Notre incapacité à montrer la capacité discriminante de la distance à la mesure par rapport à la distance de Wasserstein peut indiquer que les conditions basées sur la distance de Wasserstein ne sont pas optimales. En effet, il est possible de construire des suites de mesures (μ_n) et (ν_n) telles que la distance de Wasserstein $W_2(\mu_n, \nu_n)$ converge mais l'erreur relative entre $d_{\mu_n, m}$ et $d_{\nu_n, m}$ diverge. Nous montrons par ailleurs que la présence de bruit ambiant rend parfois impossible d'exprimer des conditions d'échantillonnage en fonction de la distance de Wasserstein alors que l'utilisation de la distance à la mesure permet d'inférer, au moins partiellement, le diagramme de persistance de l'objet. Nous définissons donc de nouvelles conditions d'échantillonnage.

Définition B.8 ⁸ Soit $M \subset \mathbb{R}^d$ une sous-variété riemannienne et soit μ une mesure de probabilité. Pour un $m \in [0, 1[$ fixé, μ est un (ϵ, r) -échantillon de M si :

$$\epsilon \geq \sup_{x \in M} d_{\mu, m}(x)$$

$$r \leq \sup \{ \ell \in \mathbb{R} \mid \forall x, d_{\mu, m}(x) < \ell \implies d(x, M) \leq d_{\mu, m}(x) + \epsilon \}$$

Par extension, si $P \subset \mathbb{R}^d$ est un nuage de points, nous dirons que P est un (ϵ, r) -échantillon de M si la mesure empirique sur P est un (ϵ, r) -échantillon de M .

Ces conditions ont la même structure que celles données en utilisant la distance de Hausdorff lorsque la distance au nuage de point est utilisée. Nous montrons que la plupart des conditions d'échantillonnage standard induisent une expression sous cette forme et donnons les valeurs des paramètres. Définissant la notion de filtration δ -effondrée $\tilde{\mathcal{F}}$ qui correspond à une troncature de la filtration \mathcal{F} au niveau du paramètre δ , c'est-à-dire le remplacement, pour $\alpha > \delta$, de F_α par l'espace ambiant, nous pouvons exprimer une inférence partielle des diagrammes de persistance.

Théorème B.9 ⁹ Soit μ un (ϵ, r) -échantillon d'une sous-variété riemannienne compacte $M \subset$

⁷Théorème 5.14

⁸Définition 6.1

⁹Théorème 6.10

\mathbb{R}^d . Les filtrations $(r + \epsilon)$ -effondrées $\tilde{\mathcal{F}} = \{\tilde{F}_\alpha\}$ de $\{d_M^{-1}([0, \alpha])\}$ et $\tilde{\mathcal{G}} = \{\tilde{G}_\alpha\}$ de $\{d_{\mu, m}^{-1}([0, \alpha])\}$ sont ϵ -entrelacées et donc :

$$d_b(\text{Dgm}(\tilde{\mathcal{F}}), \text{Dgm}(\tilde{\mathcal{G}})) \leq \epsilon.$$

B.5.4 Analyse de champs scalaires et données incomplètes

L'homologie persistante a été utilisée pour analyser des champs scalaires [32]. Il s'agit d'étudier la structure d'une fonction f à valeurs réelles définie sur une variété riemannienne M à partir d'un échantillonnage fini. Les travaux précédents ont proposé un algorithme réalisant une analyse robuste dans le cas d'un bruit sur les points borné ou d'une erreur bornée sur les valeurs fonctionnelles. Nous étendons le champ d'application de l'algorithme aux bruits aberrants en appliquant un prétraitement.

Le bruit géométrique est traité en utilisant la distance à la mesure. Les nouvelles conditions d'échantillonnage permettent de garantir une élimination des points aberrants géométriquement et une inférence correcte du support M de la fonction f . La fonction f est estimée à l'aide d'une méthode construite sur la distance à la mesure. Cette méthode considère les k plus proches voisins d'un point et recherche parmi ces points, le sous-ensemble de k' points dont les valeurs fonctionnelles présentent la plus petite variance. La valeur estimée est alors la moyenne de ces k' valeurs. L'algorithme obtenu produit des garanties en présence d'un bruit géométrique et fonctionnel contenant des valeurs aberrantes.

Théorème B.10 ¹⁰ Soient M une variété riemannienne et f une fonction c -Lipschitz sur M . Soient P un (ϵ, r) -échantillon de M et \tilde{f} des valeurs fonctionnelles tels que (7.5) est satisfaite, où $\eta \geq \epsilon$. L'algorithme traitant le bruit combiné a les garanties suivantes :
 Pour tout $\delta \in \left[2\eta + 6\epsilon, \frac{\varrho(M)}{2}\right]$ et tout $\delta' \in \left[2\eta + 2\epsilon + \frac{8}{3} \frac{r_M}{r_M - (\eta + \epsilon)} \delta, \frac{3}{4} \frac{r_M - (\eta + \epsilon)}{r_M} \varrho(M)\right]$, $H_*(f)$ et $H_*(R_\delta(\hat{L}_\alpha) \hookrightarrow R_{\delta'}(\hat{L}_\alpha))$ sont $\left(\frac{4}{3} \frac{cr_M\delta'}{r_M - (\eta + \epsilon)} + \xi s\right)$ -entrelacés où $\xi = 1$ en utilisant la médiane et $\xi = \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right)$ en utilisant la nouvelle méthode.

Bien que n'ayant pas de meilleures garanties que certaines méthodes existantes, l'utilisation de la médiane des valeurs des k plus proches voisins par exemple, notre nouvelle méthode donnent de meilleurs résultats dans certains cas intéressants, en particulier lorsque la constante de Lipschitz est grande, ou que la fonction n'est pas Lipschitz. Cela s'illustre notamment dans le débruitage d'images.

Ces résultats sont donnés avec des vitesses de convergences pour notre nouvel estimateur. Ce dernier permet également d'envisager la reconstruction de jeux de données dont une partie est manquante et en se plaçant dans des conditions non traitées par les méthodes existantes. En particulier, nous pouvons considérer des fonctions non linéaires et n'avons pas besoin de conditions sur la dimension de l'espace couvert par les valeurs.

B.6 Organisation de la thèse

Les travaux présentés dans cette thèse sont en partie le résultat de deux collaborations donnant lieu à des publications à venir. Les chapitres 2 et 3 introduisent des notions classiques en

¹⁰Théorème 7.14

analyse topologique des données et adaptent simplement les résultats de stabilité de la distance à la mesure au cas des espaces métriques quelconques. L'exception est la partie 3.5 qui constitue un travail original.

Les chapitres 4 et 5 sont le fruit de la collaboration avec F. Chazal, S. Oudot et D. Sheehy, publiés prochainement [18], et s'intéressent à l'approximation de la distance à la mesure d'une part et au calcul du diagramme de persistance des distances de puissance d'autre part. L'approximation de la distance à la mesure est obtenue par une méthode originale. La structure de donnée pour le calcul du diagramme de persistance est une adaptation technique des travaux précédents de D. Sheehy [88].

Le chapitre 7 analysant les champs scalaires est le fruit d'une visite auprès de T. Dey et Y. Wang à The Ohio State University. Les résultats obtenus en collaborant avec F. Chazal, F. Fan et S. Oudot constituent une adaptation de [32] à de nouvelles conditions de bruit par l'introduction d'une estimation originale des valeurs fonctionnelles. Les résultats sont à paraître [17].

Enfin, les chapitres 6 et 8 présentent des travaux connexes analysant les conditions d'échantillonnage utilisées au chapitre 7 et les propriétés de l'estimateur introduit dans ce même chapitre. Le chapitre 8 présente également une ouverture en direction du traitement des données incomplètes.

Bibliography

- [1] <https://www.data.gouv.fr>.
- [2] Directive 2003/98/ce on the re-use of public sector information, 2003.
- [3] Open government act of 2007, 2007.
- [4] Aaron Adcock, Daniel Rubin, and Gunnar Carlsson. Classification of hepatic lesions using the matching metric. *Computer Vision and Image Understanding*, 121:36–42, 2014.
- [5] Dana Angluin and Leslie G Valiant. Fast probabilistic algorithms for hamiltonian circuits and matchings. In *Proceedings of the ninth annual ACM symposium on Theory of computing*, pages 30–41. ACM, 1977.
- [6] Dominique Attali, André Lieutier, and David Salinas. Efficient data structure for representing and simplifying simplicial complexes in high dimensions. *International Journal of Computational Geometry & Applications*, 22(04):279–303, 2012.
- [7] Franz Aurenhammer and Hiroshi Imai. Geometric relations among voronoi diagrams. *Geometriae Dedicata*, 27(1):65–75, 1988.
- [8] Ulrich Bauer, Michael Kerber, and Jan Reininghaus. Clear and compress: Computing persistent homology in chunks. *arXiv preprint arXiv:1303.0477*, 2013.
- [9] Paul Bendich, Sang Chin, Jesse Clarke, Jonathan deSena, John Harer, Elizabeth Munch, Andrew Newman, David Porter, David Rouse, Nate Strawn, et al. Topological and statistical behavior classifiers for tracking applications. *arXiv preprint arXiv:1406.0214*, 2014.
- [10] Jesse Berwald, Marian Gidea, and Mikael Vejdemo-Johansson. Automatic recognition and tagging of topologically different regimes in dynamical systems. *arXiv preprint arXiv:1312.2482*, 2013.
- [11] Jean-Daniel Boissonnat and Frédéric Cazals. Smooth surface reconstruction via natural neighbour interpolation of distance functions. In *Proceedings of the sixteenth annual symposium on Computational geometry*, pages 223–232. ACM, 2000.

- [12] Jean-Daniel Boissonnat, Tamal K Dey, and Clément Maria. The compressed annotation matrix: An efficient data structure for computing persistent cohomology. In *Algorithms-ESA 2013*, pages 695–706. Springer, 2013.
- [13] Jean-Daniel Boissonnat, Leonidas J Guibas, and Steve Y Oudot. Manifold reconstruction in arbitrary dimensions using witness complexes. *Discrete & Computational Geometry*, 42(1):37–70, 2009.
- [14] Jean-Daniel Boissonnat, Clément Maria, et al. Computing persistent homology with various coefficient fields in a single pass. In *European Symposium on Algorithms*, 2014.
- [15] Jean-Daniel Boissonnat and Mariette Yvinec. *Algorithmic geometry*, 1998. Cambridge University Press, Cambridge.
- [16] Peter Bubenik. Statistical topology using persistence landscapes. *arXiv preprint arXiv:1207.6437*, 2012.
- [17] Mickaël Buchet, Frédéric Chazal, Tamal Krishna Dey, Fengtao Fan, Steve Oudot, and Yusu Wang. Topological analysis of scalar fields with outliers.
- [18] Mickaël Buchet, Frédéric Chazal, Steve Oudot, and Donald R Sheehy. Efficient and robust persistent homology for measures. In *Proceedings of the 26th ACM-SIAM symposium on Discrete algorithms*. SIAM, 2015.
- [19] Gunnar Carlsson and Vin De Silva. Zigzag persistence. *Foundations of computational mathematics*, 10(4):367–405, 2010.
- [20] Gunnar Carlsson, Vin de Silva, and Dmitriy Morozov. Zigzag persistent homology and real-valued functions. In *Proceedings of the twenty-fifth annual symposium on Computational geometry*, pages 247–256. ACM, 2009.
- [21] Frédéric Cazals and Joachim Giesen. Delaunay triangulation based surface reconstruction. In *Effective Computational Geometry for Curves and Surfaces*, pages 231–276. Springer, 2006.
- [22] Frédéric Chazal, David Cohen-Steiner, Marc Glisse, Leonidas J Guibas, and Steve Y Oudot. Proximity of persistence modules and their diagrams. In *Proceedings of the twenty-fifth annual symposium on Computational geometry*, pages 237–246. ACM, 2009.
- [23] Frédéric Chazal, David Cohen-Steiner, Leonidas J Guibas, Facundo Mémoli, and Steve Y Oudot. Gromov-hausdorff stable signatures for shapes using persistence. In *Computer Graphics Forum*, volume 28, pages 1393–1403. Wiley Online Library, 2009.
- [24] Frédéric Chazal, David Cohen-Steiner, and André Lieutier. A sampling theory for compact sets in euclidean space. *Discrete & Computational Geometry*, 41(3):461–479, 2009.
- [25] Frédéric Chazal, David Cohen-Steiner, and Quentin Mérigot. Geometric inference for probability measures. *Foundations of Computational Mathematics*, 11(6):733–751, 2011.

- [26] Frédéric Chazal, Vin de Silva, Marc Glisse, and Steve Oudot. The structure and stability of persistence modules. *arXiv preprint arXiv:1207.3674*, 2012.
- [27] Frédéric Chazal, Vin De Silva, and Steve Oudot. Persistence stability for geometric complexes. *Geometriae Dedicata*, pages 1–22, 2013.
- [28] Frédéric Chazal, Brittany Terese Fasy, Fabrizio Lecci, Bertrand Michel, Alessandro Rinaldo, and Larry Wasserman. Subsampling methods for persistent homology. *arXiv preprint arXiv:1406.1901*, 2014.
- [29] Frédéric Chazal, Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, Aarti Singh, and Larry Wasserman. On the bootstrap for persistence diagrams and landscapes. *arXiv preprint arXiv:1311.0376*, 2013.
- [30] Frédéric Chazal, Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, and Larry Wasserman. Stochastic convergence of persistence landscapes and silhouettes. *arXiv preprint arXiv:1312.0308*, 2013.
- [31] Frédéric Chazal, Marc Glisse, Catherine Labruère, and Bertrand Michel. Convergence rates for persistence diagram estimation in topological data analysis. In *Proceedings of The 31st International Conference on Machine Learning*, pages 163–171, 2014.
- [32] Frédéric Chazal, Leonidas J Guibas, Steve Y Oudot, and Primoz Skraba. Scalar field analysis over point cloud data. *Discrete & Computational Geometry*, 46(4):743–775, 2011.
- [33] Frédéric Chazal, Leonidas J Guibas, Steve Y Oudot, and Primoz Skraba. Persistence-based clustering in riemannian manifolds. *Journal of the ACM (JACM)*, 60(6):41, 2013.
- [34] Frédéric Chazal and André Lieutier. Topology guaranteeing manifold reconstruction using distance function to noisy data. In *Proceedings of the twenty-second annual symposium on Computational geometry*, pages 112–118. ACM, 2006.
- [35] Frédéric Chazal and Steve Yann Oudot. Towards persistence-based reconstruction in euclidean spaces. In *Proceedings of the twenty-fourth annual symposium on Computational geometry*, pages 232–241. ACM, 2008.
- [36] Bernard Chazelle. *The discrepancy method: randomness and complexity*. Cambridge University Press, 2000.
- [37] Siu-Wing Cheng, Tamal K Dey, and Edgar A Ramos. Manifold reconstruction from point samples. In *Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1018–1027. Society for Industrial and Applied Mathematics, 2005.
- [38] Moo K Chung, Peter Bubenik, and Peter T Kim. Persistence diagrams of cortical surface data. In *Information Processing in Medical Imaging*, pages 386–397. Springer, 2009.
- [39] Kenneth L Clarkson and Peter W Shor. Applications of random sampling in computational geometry, ii. *Discrete & Computational Geometry*, 4(1):387–421, 1989.

- [40] David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams. *Discrete & Computational Geometry*, 37(1):103–120, 2007.
- [41] Dorin Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 24(5):603–619, 2002.
- [42] William Crawley-Boevey. Decomposition of pointwise finite-dimensional persistence modules. *arXiv preprint arXiv:1210.0819*, 2012.
- [43] Tran Kai Frank Da, Sébastien Lorient, and Mariette Yvinec. 3D alpha shapes. In *CGAL User and Reference Manual*. CGAL Editorial Board, 4.2 edition, 2013. http://www.cgal.org/Manual/4.2/doc_html/cgal_manual/packages.html#Pkg:AlphaShapes3.
- [44] Vin De Silva, Dmitriy Morozov, and Mikael Vejdemo-Johansson. Dualities in persistent (co) homology. *Inverse Problems*, 27(12):124003, 2011.
- [45] Vin De Silva, Dmitriy Morozov, and Mikael Vejdemo-Johansson. Persistent cohomology and circular coordinates. *Discrete & Computational Geometry*, 45(4):737–759, 2011.
- [46] Mary-Lee Dequeant, Sebastian Ahnert, Herbert Edelsbrunner, Thomas MA Fink, Earl F Glynn, Gaye Hattem, Andrzej Kudlicki, Yuriy Mileyko, Jason Morton, Arcady R Mushegian, et al. Comparison of pattern detection methods in microarray time series of the segmentation clock. *PLoS One*, 3(8):e2856, 2008.
- [47] Tamal K Dey. Curve and surface reconstruction: Algorithms with mathematical analysis. 2011.
- [48] Tamal K Dey, Fengtao Fan, and Yusu Wang. Computing topological persistence for simplicial maps. *arXiv preprint arXiv:1208.5018*, 2012.
- [49] Tamal K Dey and Kuiyu Li. Topology from data via geodesic complexes. Technical report, Tech Report OSU-CISRC-3/09-TR05, 2009.
- [50] Tamal K Dey, Jian Sun, and Yusu Wang. Approximating cycles in a shortest basis of the first homology group from point data. *Inverse Problems*, 27(12):124004, 2011.
- [51] Tamal Krishna Dey, Fengtao Fan, and Yusu Wang. Graph induced complex on point data. In *Proceedings of the 29th annual symposium on Symposium on computational geometry*, pages 107–116. ACM, 2013.
- [52] Manfredo P Do Carmo. *Riemannian geometry*. Springer, 1992.
- [53] V Dobrić and JE Yukich. Asymptotics for transportation cost in high dimensions. *Journal of Theoretical Probability*, 8(1):97–118, 1995.
- [54] Yiqiu Dong and Shufang Xu. A new directional weighted median filter for removal of random-valued impulse noise. *Signal Processing Letters, IEEE*, 14(3):193–196, 2007.

-
- [55] RM Dudley. The speed of mean glivenko-cantelli convergence. *The Annals of Mathematical Statistics*, pages 40–50, 1969.
 - [56] Herbert Edelsbrunner. The union of balls and its dual shape. In *Proceedings of the ninth annual symposium on Computational geometry*, pages 218–231. ACM, 1993.
 - [57] Herbert Edelsbrunner and John Harer. *Computational topology: an introduction*. American Mathematical Soc., 2010.
 - [58] Kevin J Emmett and Raul Rabadan. Characterizing scales of genetic recombination and antibiotic resistance in pathogenic bacteria using topological data analysis. *arXiv preprint arXiv:1406.1219*, 2014.
 - [59] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Kdd*, volume 96, pages 226–231, 1996.
 - [60] Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, Larry Wasserman, Sivaraman Balakrishnan, and Aarti Singh. Statistical inference for persistent homology: Confidence sets for persistence diagrams. *arXiv preprint arXiv:1303.7117*, 2013.
 - [61] Herbert Federer. Curvature measures. *Transactions of the American Mathematical Society*, pages 418–491, 1959.
 - [62] Imola K Fodor. A survey of dimension reduction techniques, 2002.
 - [63] Sylvestre Gallot, Dominique Hulin, and Jacques Lafontaine. *Riemannian geometry*, volume 3. Springer, 1990.
 - [64] Jennifer Gamble and Giseon Heo. Exploring uses of persistent homology for statistical analysis of landmark-based shape data. *Journal of Multivariate Analysis*, 101(9):2184–2199, 2010.
 - [65] Marcio Gameiro, Yasuaki Hiraoka, Shunsuke Izumi, Miroslav Kramar, Konstantin Mischaikow, and Vidit Nanda. A topological measurement of protein compressibility. *preprint*, 2013.
 - [66] Alfred Gray et al. The volume of a small geodesic ball of a riemannian manifold. *The Michigan Mathematical Journal*, 20(4):329–344, 1974.
 - [67] Leonidas Guibas, Dmitriy Morozov, and Quentin Mérigot. Witnessed k-distance. *Discrete & Computational Geometry*, 49(1):22–45, 2013.
 - [68] László Györfi and Adam Krzyżak. *A distribution-free theory of nonparametric regression*. Springer, 2002.
 - [69] Allen Hatcher. *Algebraic topology*. Cambridge University Press, 2002.

- [70] Giseon Heo, Jennifer Gamble, and Peter T Kim. Topological analysis of variance and the maxillary complex. *Journal of the American Statistical Association*, 107(498):477–492, 2012.
- [71] Hugues Hoppe, Tony DeRose, Tom Duchamp, John McDonald, and Werner Stuetzle. *Surface reconstruction from unorganized points*, volume 26. ACM, 1992.
- [72] Danijela Horak, Slobodan Maletić, and Milan Rajković. Persistent homology of complex networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2009(03):P03034, 2009.
- [73] Joseph Horowitz and Rajeeva L Karandikar. Mean rates of convergence of empirical measures in the wasserstein metric. *Journal of Computational and Applied Mathematics*, 55(3):261–273, 1994.
- [74] Samory Kpotufe. k-nn regression adapts to local intrinsic dimension. In *Advances in Neural Information Processing Systems*, pages 729–737, 2011.
- [75] M Kramar, A Goulet, L Kondic, and K Mischaikow. Persistence of force networks in compressed granular media. *Physical Review E*, 87(4):042207, 2013.
- [76] Hyekyoung Lee, Hyejin Kang, Moo K Chung, Bung-Nyun Kim, and Dong Soo Lee. Persistent brain network homology from the perspective of dendrogram. *Medical Imaging, IEEE Transactions on*, 31(12):2267–2277, 2012.
- [77] Chunyuan Li, Maks Ovsjanikov, and Frederic Chazal. Persistence-based structural recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1995–2002, 2013.
- [78] Stuart Lloyd. Least squares quantization in pcm. *Information Theory, IEEE Transactions on*, 28(2):129–137, 1982.
- [79] Ching-Ta Lu and Tzu-Chun Chou. Denoising of salt-and-pepper noise corrupted image using modified directional-weighted-median filter. *Pattern Recognition Letters*, 33(10):1287–1295, 2012.
- [80] Quentin Mérigot. Lower bounds for k-distance approximation. In *Proceedings of the 29th annual symposium on Symposium on computational geometry*, pages 435–440. ACM, 2013.
- [81] Nikola Milosavljević, Dmitriy Morozov, and Primož Skraba. Zigzag persistent homology in matrix multiplication time. In *Proceedings of the twenty-seventh Annual Symposium on Computational Geometry*, pages 216–225. ACM, 2011.
- [82] David M. Mount and Sunil Arya. ANN: Library for approximate nearest neighbour searching. 1998.

-
- [83] James R Munkres. *Elements of algebraic topology*, volume 2. Addison-Wesley Reading, 1984.
 - [84] Monica Nicolau, Arnold J Levine, and Gunnar Carlsson. Topology based data analysis identifies a subgroup of breast cancers with a unique mutational profile and excellent survival. *Proceedings of the National Academy of Sciences*, 108(17):7265–7270, 2011.
 - [85] Sylvain Paris and Frédo Durand. A topological approach to hierarchical segmentation using mean shift. In *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*, pages 1–8. IEEE, 2007.
 - [86] Giovanni Petri, Martina Scolamiero, Irene Donato, and Francesco Vaccarino. Topological strata of weighted complex networks. *PloS one*, 8(6):e66506, 2013.
 - [87] Walter Rudin. *Real and complex analysis (3rd)*. New York: McGraw-Hill Inc, 1986.
 - [88] Donald R Sheehy. Linear-size approximations to the vietoris–rips filtration. *Discrete & Computational Geometry*, 49(4):778–796, 2013.
 - [89] Gurjeet Singh, Facundo Memoli, Tigran Ishkhanov, Guillermo Sapiro, Gunnar Carlsson, and Dario L Ringach. Topological analysis of population activity in visual cortex. *Journal of vision*, 8(8):11, 2008.
 - [90] Primož Skraba, Maks Ovsjanikov, Frederic Chazal, and Leonidas Guibas. Persistence-based segmentation of deformable shapes. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on*, pages 45–52. IEEE, 2010.
 - [91] Thierry Sousbie. The persistent cosmic web and its filamentary structure–i. theory and implementation. *Monthly Notices of the Royal Astronomical Society*, 414(1):350–383, 2011.
 - [92] Thierry Sousbie, Christophe Pichon, and Hajime Kawahara. The persistent cosmic web and its filamentary structure–ii. illustrations. *Monthly Notices of the Royal Astronomical Society*, 414(1):384–403, 2011.
 - [93] Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
 - [94] Cédric Villani. *Topics in optimal transportation*. Number 58. American Mathematical Soc., 2003.
 - [95] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
 - [96] Shuenn-Shyang Wang and Cheng-Hao Wu. A new impulse detection and filtering method for removal of wide range impulse noises. *Pattern Recognition*, 42(9):2194–2202, 2009.

Bibliography

- [97] Yuan Wang, Hernando Ombao, and Moo K Chung. Persistence landscape of functional signal and its application to epileptic electroencephalogram data.
- [98] Hui Zhang, Jason E Fritts, and Sally A Goldman. Image segmentation evaluation: A survey of unsupervised methods. *computer vision and image understanding*, 110(2):260–280, 2008.
- [99] Afra Zomorodian and Gunnar Carlsson. Computing persistent homology. *Discrete & Computational Geometry*, 33(2):249–274, 2005.

Index

- (ϵ, r) -sample, 92
- (k, k', Δ) -functional-sample, 106
- δ -collapsed filtration, 99
- δ -matching, 16
- ϵ -correspondence, 71
- ϵ -homomorphism, 18
- σ -algebra, 27
- q -tameness, 14

- abstract simplex, 9
- abstract simplicial complex, 9

- barcode, 16
- bottleneck distance, 17
- boundary homomorphism, 11

- chain complex, 10
- clique complex, 10
- clutter noise, 94
- contiguity, 13, 79

- discrepancy, 108, 121
- distance to a measure, 27, 29
- doubling constant, 22
- doubling dimension, 22

- empirical measure, 28

- facet, 9
- filtered simplicial complex, 20
- filtration, 14

- Günther-Bishop theorem, 24
- general position, 22
- geometric complex, 10

- geometric simplex, 9

- Hausdorff distance, 21
- homology group, 11, 12

- incomplete data, 126
- interleaving, 18

- logarithmic bottleneck distance, 18

- measure, 27
- medial axis, 22
- median, 108, 123
- metric space, 21
- minimizing geodesic, 23

- persistence, 14
- persistence diagram, 15
- persistence module, 14
- power cell, 45
- power diagram, 45
- power distance, 44
- probability measure, 27
- pseudo-distance to a measure, 29

- reach, 22
- regression, 121
- Riemannian ϵ -sampling, 23

- sectional curvature, 24
- simplicial homology, 11
- simplicial map, 13
- singular homology, 13
- singular simplex, 13
- smooth submanifold, 23
- sparse Rips filtration, 75
- sparse weighted Rips filtration, 75

Index

strong convexity radius, 24

submeasure, 31

support, 28

transport plan, 28

triangulable space, 14

Voronoi cell, 41

Voronoi diagram, 41

Wasserstein distance, 29

weighted Čech complex, 67

weighted Rips complex, 68

witnessed k -distance, 51

zig-zag persistence, 20